

**The Enigma of Genetic Linkage in Molecular Breeding for Maize**

A DISSERTATION

SUBMITTED TO THE FACULTY OF THE GRADUATE SCHOOL

OF THE UNIVERSITY OF MINNESOTA

BY

Joshua Andrew Sleper

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS

FOR THE DEGREE OF

DOCTOR OF PHILOSOPHY

Rex Bernardo, Advisor

June 2017



## **Acknowledgements**

I would like to thank my advisor, Dr. Rex Bernardo, for his guidance, support, and dedication to the process of teaching and developing the next generation of plant scientists. I would also like to thank my advisory committee, Todd Warner, Drs. Jim Anderson, Yang Da, and Candice Hirsch, for their input on this project and their commitment to my development as a scientist. Also, I would like to thank my scientific mentors over the years, Todd Warner, Drs. Rex Bernardo, Craig Butler, Gary Muehlbauer, Dietrich Borchardt, who have taught me scientific thought and given great inspiration.

I am very thankful, as a friend and a colleague, for the time spent with fellow lab members and graduate students Amy Jacobson, Karl Butenhoff, Jean-Michel Michno, Lian Lian, Nick Ames, Elizabeth Blissett, and Sofia Brandariz Zerboni.

I would like to thank Syngenta who made this possible. Additionally, I am very thankful for all of my colleagues at Syngenta, Raffaele Capitanio, Jason Cromley, Giovanni Della Porta, Ben Ford, Sharon Gergen, Shreyartha Mukherjee, and Zach Toland, who have given help and guidance on this project.

Lastly, I would like to thank all of my friends, family, and church for all of the support and encouragement throughout this time. Finally, I would like to thank my wife Amy for her faithfulness and love.

## **Abstract**

Linkage among quantitative trait loci prevents the release of hidden genetic variation, but also preserves desirable gene combinations. This dissertation, which includes three studies, shows the continuing enigma of linkage in maize (*Zea mays* L.) breeding. The first study aimed to determine if the additional recombinations in doubled haploids induced from F<sub>2</sub> instead of F<sub>1</sub> plants leads to a larger genetic variance and a superior mean of the best lines. In two maize populations, inducing doubled haploids from F<sub>2</sub> plants did not improve the mean, and it increased the genetic variance for moisture, but not for yield and plant height. The second study aimed to determine if multi-allelic markers or haplotypes improve the prediction accuracy of genomewide selection in three-way breeding populations, which could have three alleles per locus. In both simulated and empirical maize populations, accounting for multiple alleles did not improve the prediction accuracy over a biallelic model. The third study aimed to determine if genomewide markers can be used to partition trait effects into independent and correlated portions, and if selection on the independent portion was more effective than selection on the entire trait. Results from four cycles of selection showed that selection only for the independent portion did not lead to higher responses for yield, moisture, and plant height. Overall, genetic linkage both assists and confounds molecular breeding efforts in maize.

## Table of Contents

<b>List of Tables .....</b>	<b>v</b>
<b>List of Figures .....</b>	<b>vii</b>
<b>Chapter 1: Recombination and genetic variance among maize doubled haploids induced from F<sub>1</sub> and F<sub>2</sub> plants .....</b>	<b>1</b>
Introduction .....	2
Materials and Methods .....	4
Germplasm and Field Experiments .....	4
Single Nucleotide Polymorphism Analysis .....	5
Data Analysis .....	6
Simulations .....	7
Results .....	8
Discussion .....	10
Expected recombination frequencies, means, and genetic variances with DHF1 and DHF2 lines .....	10
Observed versus expected recombination frequencies, means, and genetic variances .....	12
Mean of best 10% of lines .....	14
<b>Chapter 2: Genomewide Selection with Biallelic versus Triallelic Models in Three-way Maize Populations .....</b>	<b>21</b>
Introduction .....	22
Materials and Methods .....	24
Biallelic Model .....	24

Marker Interval Model .....	25
Allele Phasing Model .....	25
Simulation Experiments .....	27
Maize Empirical Data .....	29
Results and Discussion .....	31
Simulated Populations .....	31
Empirical Populations .....	31
Equivalency between the Models in Practice .....	32
Application .....	34
<b>Chapter 3: Genomewide selection for unfavorably correlated traits in maize .....</b>	<b>41</b>
Introduction .....	42
Materials and Methods .....	44
Training Data .....	44
Genomewide Selection Models .....	46
Recurrent Selection .....	47
Response to Selection .....	49
Results and Discussion .....	51
Training Data .....	51
Response to Genomewide Selection and Level of Inbreeding .....	53
Conclusions .....	55
<b>Bibliography .....</b>	<b>61</b>

## List of Tables

Table 1: Trait means, genetic variance ( $V_G$ ), residual variance ( $V_R$ ), heritability ( $h^2$ ), and mean of best 10% of testcrosses ( $U_{0.10}$ ) in DHF1 and DHF2 maize populations. .	16
Table 2: Mean and variance of the number of recombinations per chromosome and the frequency of inheriting an entire parental chromosome among maize DHF1 and DHF2 lines. ....	17
Table 3: Repulsion:coupling ratios (95% confidence intervals in parenthesis) among DHF1 and DHF2 maize lines. ....	18
Table 4: Estimated repulsion:coupling (R:C) ratios and correlation between underlying and estimated R:C ratios among simulated DHF1 and DHF2 lines. ....	19
Table 5: Estimated prediction accuracy in simulated three-way populations. ....	36
Table 6: Summary statistics for assessing the predictive ability of genomewide selection in each of four three-way maize test populations. ....	37
Table 7: Predictive ability in maize three-way test populations for yield, moisture, and test weight. ....	38
Table 8: Estimated prediction accuracy in simulated three-way population after six generations of random mating, under three genetic models that differed in heritability ( $h^2$ ) and number of quantitative trait loci (QTL). ....	39
Table 9: Summary statistics for training populations used for genomewide selection in maize. ....	57
Table 10: Pairwise phenotypic correlations between the full trait ( $y$ ), the independent portion of the trait ( $y_n$ ), and the correlated portion of the trait ( $y_R$ ). Genotypic correlations calculated using analysis of covariance ( $r_G$ ) and using molecular	

markers for the full trait ( $g$ ), the independent portion of the trait ( $g_n$ ), and the correlated portion of the trait ( $g_R$ ). .....	58
Table 11: Testcross performance of different cycles of genomewide selection for maize Population 1. ....	59
Table 12: Inbreeding level and mean heterozygosity between the Control Model and the Independent Model for both maize populations. ....	60



## List of Figures

- Figure 1: Two-locus models for the expected mean and genetic variance among DHF1 lines (solid line) and DHF2 lines (dashed line). The genotypic values were 2 for the *A-B-* genotype and 0 for all other genotypes with complementary epistasis; and 0 for the *aabb* genotype and 2 for all other genotypes with dominant epistasis. .... 20
- Figure 2: Allele Phasing in a three-way population: (A) inbred parents are genotyped with biallelic SNP markers; (B) inbred progeny of the three-way population are genotyped with biallelic SNP markers; (C) genotypes for a three-allele model are inferred in the three-way progeny according to allelic descent; and (D) final genotypes are projected based on flanking marker information, assuming no double recombination. .... 40

## **Chapter 1: Recombination and genetic variance among maize doubled haploids induced from F<sub>1</sub> and F<sub>2</sub> plants**

Maize (*Zea mays* L.) breeders rely on doubled haploid (DH) technology for fast and efficient production of inbred lines. Breeders can induce DH lines most quickly from F<sub>1</sub> plants (DHF1), or induce DH lines from F<sub>2</sub> plants (DHF2) to allow selection prior to DH induction and have more recombinations. Our objective was to determine if the additional recombinations in maize DHF2 lines leads to a larger genetic variance and a superior mean of the best lines. A total of 311 DHF1 and 241 DHF2 lines, derived from the same biparental cross, were crossed to two testers and evaluated in multilocation trials in Europe and the U.S. The mean number of recombinations per genome was 14.48 among the DHF1 lines and 21.38 among the DHF2 lines. The means of the DHF1 and DHF2 lines did not differ for yield, moisture, and plant height. The genetic variance was higher among DHF2 lines than among DHF1 lines for moisture, but not for yield and plant height. Additionally, the repulsion:coupling ratio (a new statistic derived from genomewide marker effects) was higher among DHF1 lines than among DHF2 lines for moisture, but not for yield and plant height. However, the higher genetic variance for moisture among DHF2 lines did not lead to a lower moisture of the best 10% of the lines. Our results indicated that the decision of whether to induce DH lines from F<sub>1</sub> or F<sub>2</sub> plants needs to be made from considerations other than the performance of the resulting DHF1 or DHF2 lines.

## Introduction

Doubled haploid (DH) technology has decreased the time needed to create a maize (*Zea mays* L.) inbred from 6-7 selfing generations to two generations. As a result, commercial maize breeding has become more efficient in terms of genetic gain per year (Longin, 2008). In a typical maize breeding program,  $F_1$  or  $F_2$  plants from a biparental cross are crossed to a haploid inducer. Haploid progeny are identified using a morphological marker and their chromosomes are doubled by treating seeds, embryos, or seedlings with colchicine (Rober et al., 2005). The resulting DH lines are then testcrossed and evaluated for their hybrid performance within three years from the time that the initial  $F_1$  is made.

Maize breeders need to decide whether to induce and create DH lines from  $F_1$  plants or from  $F_2$  plants (Bernardo, 2009). Haploid induction in either generation has both advantages and disadvantages that pertain to three factors: time needed to create DH lines; amount of recombination; and ability to select plants prior to haploid induction. First, inbred development is quickest with  $F_1$ -induced DH lines (DHF1), which are created in two life cycles ( $F_1$  and induction generations). On the other hand,  $F_2$ -induced DH lines (DHF2) are created in three life cycles, with the additional generation ( $F_1$  to  $F_2$ ) causing a slight delay in inbred development. Second, because a DHF1 line is created after a single generation of meiosis, the line is characterized by only about 10 crossover events across the genome (Smith et al., 2008). Empirical results have shown that DHF1 lines inherit one or more fully intact parental chromosomes more than 95% of the time (Smith et al., 2008). But as we show later in this article, about 50% more crossovers are expected in DHF2 lines than in DHF1 lines. The disruption of unfavorable linkages is therefore greater in DHF2 lines than in DHF1 lines, but the preservation of favorable linkages is less in DHF2 lines

than in DHF1 lines. Third, selection prior to haploid induction can be done among F<sub>2</sub> plants but not among F<sub>1</sub> plants, which are genetically identical. Maize breeders must weigh these advantages and disadvantages in choosing the generation for inducing DH lines.

The fewer recombinations in DHF1 lines than in recombinant inbreds could create a genetic bottleneck that would reduce long-term genetic gains (Jannink and Abadie, 1999). Simulation results have shown larger long-term (fifteen cycles of selection) genetic gains with DHF2 lines than with DHF1 lines (Bernardo, 2009). The simulation suggested using DHF2 lines as a potential compromise between the speed of creating DHF1 lines and the additional recombinations found in recombinant lines. In lieu of these findings, empirical data is needed to compare DHF1 and DHF2 lines in regards to the level of genetic variance ( $V_G$ ) and the amount of repulsive linkage blocks. In the short term, the additional recombinations among DHF2 lines could increase the  $V_G$  by releasing variation that is hidden by repulsion linkages (Weir et al., 1980). To illustrate, suppose loci *A* and *B* are linked in repulsion phase. The genotypic values are of 2 for *AABB*, 1 for *AAbb* and *aaBB*, and 0 for *aabb*. An *AAbb* × *aaBB* cross is made and DH lines are developed. If recombination does not occur, the DH lines will be 50% *AAbb* and 50% *aaBB*. These DH lines all have a genotypic value of 1, and no  $V_G$  is expressed across the two loci. But if, due to an additional meiosis, recombination occurs between the *A* and *B* loci, DH lines with the *AABB* and *aabb* genotypes can be recovered. The DHF2 lines have unequal genotypic values, thereby leading to a positive  $V_G$  across the two loci.

To date, there have been no published empirical studies comparing the  $V_G$  and recombination frequencies between DHF1 and DHF2 lines. An increase in  $V_G$  due to an extra generation of random mating within a population would give a breeder more incentive

to develop DHF2 lines. Our objective was to determine if the additional recombinations in maize DHF2 lines leads to a larger genetic variance and a superior mean of the best lines.

## **Materials and methods**

### *Germplasm and field experiments*

We studied a maize biparental population from Syngenta Seeds, LLC. The population was derived from two Lancaster SureCrop fully inbred parents and initially included 360 DHF1 lines and 253 DHF2 lines. Haploids were induced by crossing 150  $F_1$  plants and 150 random  $F_2$  plants to a Syngenta haploid inducer, and the chromosomes were doubled using a colchicine treatment. To avoid any selection, each of the 150 random  $F_2$  plants was crossed to the haploid inducer. Some DHF2 lines were derived from the same  $F_2$  plant, but records were not kept of the exact numbers of DHF2 lines derived from each of the 150  $F_2$  plants. Simulation experiments we conducted indicated only a 3% difference in  $V_G$  (results not shown) when two DHF2 lines were derived from an  $F_2$  plant versus when one DHF2 line was derived from an  $F_2$  plant.

The population was crossed to two different B73-type testers. The days to relative maturity were 112 for Parent 1, 116 for the Parent 2, 114 for Tester A, and 115 for Tester B. Simple matching coefficients based on a proprietary set of Syngenta single nucleotide polymorphism (SNP) markers were as follows: 0.75 between Parent 1 and Parent 2; 0.60 between Parent 1 and Tester A; 0.54 between Parent 1 and Tester B; 0.60 between Parent 2 and Tester A; 0.53 between Parent 2 and Tester B; and 0.73 between Tester A and Tester B.

Population A, which referred to the testcrosses to Tester A, was evaluated at six locations in Italy (Moscazzano, Torre di Mosto, and Casale Monferrato) and Spain (Palma del Río, Posadas, and Arroyo de San Serván) in 2013. Population B, which referred to the testcrosses to Tester B, was evaluated at seven U.S. Midwest locations (Washington, Indiana; Gold, Illinois; Dayton, Cass, Palestine, and Eagle Grove, Iowa; and Foster, Nebraska) in 2014. The testcrosses were grown along with four check hybrids in a single replication at each location. The entries were grown in two-row plots, each row 6.1 m long and spaced 76 cm apart, at a plant population density of 89,000 plants per hectare. Plant height (cm) was measured as the distance from the soil surface to the node of the flag leaf of one representative plant per plot. Yield ( $\text{Mg ha}^{-1}$  at  $155 \text{ g H}_2\text{O kg}^{-1}$ ) and moisture ( $\text{g kg}^{-1}$ ) were obtained. Yield and moisture were measured at each location while plant height was measured at six locations for Population A and three locations for Population B.

#### *Single nucleotide polymorphism analysis*

Seedling DNA was extracted from leaf punches of one plant per DH line. Each DH line was genotyped with 3072 SNP markers (proprietary to Syngenta) on the Illumina GoldenGate platform. A total of 907 markers were polymorphic between Parent 1 and Parent 2. A SNP marker was excluded if it had more than 10% missing data, 10% heterozygosity, or a minor allele frequency less than 5%. A DH line was excluded if it had more than 20% missing data or more than 10% SNPs segregating at markers monomorphic between the two parents. These criteria led to 311 DHF1 lines, 241 DHF2 lines, and 725 segregating SNP markers being used in the final analysis. A linkage map was constructed via JoinMap v. 3.0 (Ooijen and Voorrips, 2002) with the DHF1 and DHF2 lines and 725

SNP markers. The linkage map was 1782 centiMorgans and had a mean marker spacing of 2.45 centiMorgans between adjacent markers.

#### *Data analysis*

Population A and Population B were analyzed separately because they were evaluated in different environments. The linear model included the grand mean, location effect, line effect, and residual effect with each effect assumed as random. Testcross  $V_G$  and nongenetic variance ( $V_R$ ) were calculated from across-locations analysis using the “lmer” function in the “lme4” package (Bates et al., 2013). Genotype-by-environment interaction variance and within-location error variance were confounded in  $V_R$ . Entry-mean heritability was estimated separately for the DHF1 and DHF2 testcrosses as  $h^2 = V_G/(V_G + V_R/l)$ , where  $l$  was the number of locations for the trait. A t-test was used to test for significance ( $P = 0.05$ ) of the difference between the means of the DHF1 and DHF2 lines. For each trait, the observed mean of the best 10% of the lines (i.e., observed  $U_{0.10}$ ) was calculated among the DHF1 and DHF2 lines. The difference in the observed  $U_{0.10}$  among DHF1 and DHF2 lines was tested for significance via a t-test. The predicted mean of the best 10% of the lines (i.e., predicted  $U_{0.10}$ ) was also calculated from the estimates of the mean,  $V_G$ , and  $h^2$  (Melchinger et al., 1988). A 95% bootstrap confidence interval was constructed for the difference in  $V_G$  between the DHF1 and DHF2 lines. Bootstrapping was conducted by resampling the DHF1 and DHF2 lines.

The incidence of unintended selection among  $F_2$  plants was analyzed by testing the difference in SNP allele frequencies between the DHF1 and DHF2 lines. Two-sided  $z$ -tests were conducted using the “prop.test” function in R. A Bonferroni correction was applied ( $P = 0.05/725$ ) to account for multiple testing with 725 SNP loci.

Recombination events were counted for each DH line by identifying the crossover locations on each chromosome. Within each chromosome, the difference between the mean number of crossovers between DHF1 and DHF2 lines was tested for significance ( $P = 0.05$ ) via the “poisson.test” function in R.

The relative proportion of repulsion versus coupling linkages (denoted by R:C) was estimated. Genomewide marker effects were first calculated by ridge regression-best linear unbiased prediction with the rrBLUP package (Endelman, 2011). Given that the genomewide marker effects referred to those for the marker alleles from Parent 1, coupling linkage was declared when the effects of two adjacent markers were both positive or were both negative. Repulsion linkage was declared when the effects of adjacent markers had opposite signs. The number of coupling and repulsion linkages within each chromosome was calculated, and R:C was obtained as the ratio between the total number (across chromosomes) of coupling linkages and total number of repulsion linkages. The R:C values were obtained separately for the DHF1 and DHF2 lines in each population. The significance of each R:C value was tested via 95% bootstrap confidence intervals.

### *Simulations*

We conducted simulation experiments to determine the correspondence between the R:C ratios estimated from genomewide marker effects and the known R:C ratios that arise from the underlying QTL. We considered a 200 cM chromosome with 200 evenly spaced markers that were polymorphic between the two parents. The trait was controlled by 10, 19, and 28 QTL that were randomly distributed across the chromosome. The number of possible linkages was 9 with 10 QTL, 18 with 19 QTL, and 27 with 28 QTL. The QTL effects followed a geometric series (Lande and Thompson, 1990) and the QTL did not



exhibit dominance or epistasis. Different known R:C ratios were simulated: 0:9 or 0%; 2:7 or 0.28; 4:5 or 0.80, 6:3 or 2.0; and 8:1 or 8.0. These known R:C values were simulated by assigning different genotypes (out of several possible combinations) to the two parents. For example, with 10 QTL and an R:C ratio of 4:5, a possible pair of parental genotypes was *AAAbbCCDDEEffGGHHIIJJ* for the first parent and *aaBBccddeeffGghhiiijj* for the second parent.

A total of 250 DHF1 lines and 250 DHF2 lines were simulated. Nongenetic values were added to the genetic values to obtain phenotypic values. These nongenetic values had a mean of zero and a variance that was scaled to achieve  $h^2$  values of 0.30, 0.50, or 0.80. Each combination of the number of QTL, R:C level, and  $h^2$  was simulated 100 times, with the QTL locations and parental genotypes (for the same R:C ratio) differing in each repeat. The mean R:C ratio was calculated across repeats.

## Results

For both Population A and Population B, none of the differences between trait means of the DHF1 and DHF2 lines was significant ( $P = 0.05$ ). For yield, moisture, and plant height, the differences between the DHF1 and DHF2 means were less than 1% for both populations (Table 1). Additionally, the difference between the observed means of the best 10% of the DHF1 lines and best 10% of the DHF2 lines was not significant ( $P = 0.05$ ) for any of the traits.

All estimates of  $V_G$  were significant ( $P = 0.05$ ; Table 1). In both Population A and Population B,  $V_G$  for moisture was significantly ( $P = 0.05$ ) higher among DHF2 lines than among DHF1 lines (Table 1). In particular, the  $V_G$  estimate for moisture was 58% higher among DHF2 lines than among DHF1 lines in Population A, and 34% higher among DHF2

lines than among DHF1 lines in Population B. In contrast,  $V_G$  for yield and for plant height did not differ significantly between the DHF1 and DHF2 lines in both populations. The  $h^2$  for yield (0.40 to 0.54) was lower than the  $h^2$  for moisture and plant height. Within each trait and population, the DHF1 lines and DHF2 lines had similar estimates of  $V_R$ . The larger  $h^2$  for moisture among DHF2 lines than among DHF1 lines was due to differences in  $V_G$  rather than in  $V_R$ .

Only 11 out of the 725 SNP markers had a significant difference in allele frequency between the DHF1 and DHF2 lines. These significant differences in SNP allele frequencies ranged from 0.16 to 0.23. The 11 SNPs were distributed as follows (cM positions in parentheses): one on chromosome 1 and 9, two on chromosome 3 (158 and 162 cM), three on chromosome 5 (19, 187 and 194 cM), and four on chromosome 6 (1, 34, 36, and 50 cM). The mean number of recombinations across the genome was 14.48 per DHF1 line and 21.38 per DHF2 line (Table 2). The difference in the mean number of recombinations between the DHF2 and DHF1 lines was significant for each of the 10 chromosomes. The mean number of recombinations differed among the chromosomes and ranged from 0.64 in DHF1 lines and 0.89 in DHF2 lines for chromosome 2, to 2.56 in DHF1 lines and 3.53 in DHF2 lines on chromosome 5. The range in the number of recombinations per chromosome was larger in the DHF2 lines than in the DHF1 lines. Furthermore, the mean frequency of an entire chromosome being inherited was 0.23 in DHF1 lines and 0.14 in DHF2 lines (Table 2). The frequencies of an entire chromosome being inherited were lowest for chromosomes 1 and 5 (0.04 to 0.08) and highest for chromosome 2 (0.38 to 0.48).

In both Population A and Population B, the ratio of the frequencies of repulsion and coupling linkages (R:C) for moisture was significantly larger ( $P = 0.05$ ) among DHF1 lines than among DHF2 lines (Table 3). In contrast, R:C for yield and for plant height did not differ significantly between the DHF1 and DHF2 lines in both populations.

Simulations showed that the underlying R:C values were highly correlated with the R:C ratios estimated from genomewide marker effects, with the correlations ranging from 0.61 to 0.98 across different genetic models (Table 4). Compared with the underlying R:C values, which ranged from 0 to 8, the estimated R:C ratios were greatly shrunken towards zero and ranged from about 0.10 to about 0.25. The estimated R:C ratios generally increased as the number of QTL decreased and the  $h^2$  increased. For each genetic model, the estimated R:C ratios were higher (by 4 to 29%) among DHF2 lines than among DHF1 lines.

## Discussion

### *Expected recombination frequencies, means, and genetic variances with DHF1 and DHF2 lines*

The fundamental difference between DHF1 and DHF2 lines is that an additional meiotic event among  $F_2$  individuals results in additional recombinations. Because the frequency of recombinant gametes produced by  $F_1$  individuals of a dihybrid cross is  $r$ , the expected frequency of recombinants among DHF1 lines is also  $r$ . Among recombinant inbreds, the expected frequency of recombinants was previously shown to be  $2r/(1 + 2r)$  (Haldane and Waddington, 1931). By considering the frequencies and gametic outputs of the different genotypes in the  $F_2$ , we found that the frequency of recombinants among

DHF2 lines is  $r(1.5 - r)$ . As such, the ratio between the number of recombinations in DHF2 lines versus DHF1 lines is  $1.5 - r$ , whereas the ratio between the number of recombinations in recombinant inbreds versus DHF1 lines is  $2/(1 + 2r)$ . As  $r$  approaches zero, we expect recombinant inbreds to have 100% more recombinations than DHF1 lines and DHF2 lines to have 50% more recombinations than DHF1 lines. The difference in the frequency of recombinants between DHF1 lines and recombinant inbreds has been previously demonstrated by Riggs and Snape (1997) to affect the means and variances of the lines. Overall, we expected the differences between DHF1 lines and recombinant inbreds shown by Riggs and Snape (1977) to hold between DHF1 and DHF2 lines, but to a lesser degree due to the smaller difference in the number of meiotic events.

When both linkage and epistasis are present, differences in the frequencies of recombinants lead to differences in the expected means of the DHF1 and DHF2 lines (Fig. 1). Linkage alone or epistasis alone does not lead to a difference in the expected means of the DHF1 and DHF2 lines. The combinations of linkage phase, type of epistasis, and generation of DH induction that lead to a higher mean are as follows (Fig. 1): (1) coupling linkage, complementary gene action, DHF1; (2) repulsion linkage, duplicate dominant epistasis, DHF1; (3) repulsion linkage, complementary gene action, DHF2; and (4) coupling linkage, duplicate dominant epistasis, DHF2. The difference between the expected means of the DHF1 and DHF2 lines was largest with  $r = 0.25$  (Fig. 1).

Whereas both linkage and epistasis are required for the DHF1 and DHF2 lines to differ in their expected means, linkage alone causes the expected  $V_G$  to differ between the DHF1 and DHF2 lines. Coupling linkage leads to a larger  $V_G$  among DHF1 lines than among DHF2 lines for both additive and epistatic genetic models (Fig. 1). Repulsion

linkage leads to a larger  $V_G$  among DHF2 lines than among DHF1 lines regardless of the type of gene action. This higher  $V_G$  among DHF2 than among DHF1 lines corresponds to the hidden genetic variance that is released upon the disruption of repulsion linkages (Weir et al., 1980). The type of epistasis (complementary gene action or duplicate dominant epistasis) does not affect whether  $V_G$  is greater in DHF1 lines or in DHF2 lines. However, differences in the expected  $V_G$  were largest when the recombination frequency is  $r = 0.25$ . Lastly, when the two loci are unlinked, we expect equal  $V_G$  between DHF1 and DHF2 lines regardless of the genetic model.

*Observed versus expected recombination frequencies, means, and genetic variances*

The observed frequencies of recombinants and population means of the DHF1 versus DHF2 lines agreed with the expectations. The observed ratio of recombinations in the DHF2 versus DHF1 was  $21.38/14.48 = 1.477$  (Table 2). The mean distance between adjacent markers was 2.45 cM, which corresponded to  $r = 0.0244$  with the Kosambi mapping function. The observed ratio of recombinations of 1.477 was therefore very close to the expected ratio of  $(1.5 - r) = 1.475$  for DHF2 versus DHF1 lines. This result was unsurprising because the same data set was used to create the linkage map and to count the frequency of recombinants.

The lack of a significant difference between the means of DHF1 and DHF2 lines indicated that epistasis, as well as selection, was negligible. The overall lack of selection was further shown by the lack of a significant difference in allele frequencies between DHF1 and DHF2 plants at 714 out of the 725 SNP markers. The 11 SNP markers that showed significant differences in allele frequencies between the DHF1 and DHF2 lines could imply segregation distortion. Chromosomes 3, 5, and 6 each possessed at least two

significant markers, but these few regions did not seem to correspond with previously identified areas of segregation distortion (Lu et al., 2002). The overall lack of evidence for selection indicated that any observed differences in both the mean and  $V_G$  was not due to altered allele frequencies.

The observed differences in  $V_G$  among DHF1 and DHF2 lines were inconsistent across the three traits we studied. Moisture had a significantly higher  $V_G$  among DHF2 lines than among DHF1 lines; in contrast,  $V_G$  did not differ significantly between DH generations for both yield and plant height. This trend held across both testers and testing environments. Our expected results showed that DHF2 lines are expected to have a higher  $V_G$  if the quantitative trait loci (QTL) are linked in repulsion phase (Fig. 1). The observed results for  $V_G$  therefore suggested that QTL for moisture are predominantly linked in repulsion phase in the cross we studied. On the other hand, we surmise that neither linkage phase was predominant for yield and plant height in the cross that we studied. We speculate that these QTL for moisture are largely non-epistatic, as supported by the equal means of the DHF1 and DHF2 lines (Table 1) and negligible epistasis reported in maize for quantitative traits (Stuber and Moll, 1971; Silva and Hallauer, 1975; Melchinger et al., 1986).

The disruption of repulsion linkages, and the subsequent release of hidden  $V_G$ , was consistent with the R:C ratio for moisture being significantly lower for DHF2 lines than for DHF1 lines. For yield and plant height, the lack of a significant difference in  $V_G$  among DHF1 versus DHF2 lines was consistent with the lack of a significant difference in R:C ratios for these two traits. The simulation results confirmed that genomewide marker effects can be used to assess the relative frequencies between populations, but not the

absolute frequencies in each population, of repulsion and coupling linkages. The simulation also demonstrated that the number of recombinations needed to break repulsion linkages is a function of the genome size and the number of QTL affecting the trait(s) of interest (Bernardo, 2009). As the number of QTL increased with a fixed chromosome length, R:C decreased demonstrating a lower resolution in detecting repulsive elements at a high QTL density. Interpretation of an R:C ratio by itself should not be done because two SNP markers may seem to be in coupling linkage when, in fact, their effects are due to their linkage to the same QTL. For instance, suppose four markers and one QTL are in coupling phase in the following map order:  $M_1$ - $M_2$ -QTL- $M_3$ - $M_4$ . If all four markers are closely linked, the alleles inherited from the same parent will have positive effects. In calculating the R:C ratio, a total of three coupling linkages will be counted even though there is only a single QTL.

#### *Mean of best 10% of lines*

The mean of the best 10% of the lines ( $U_{0.10}$ , Table 1) combines information on the population mean,  $V_G$ , and  $h^2$  in a way that is most meaningful to a breeder. For yield and plant height, the equal means and equal  $V_G$  estimates between the DHF1 and DHF2 lines indicated that no differences in  $U_{0.10}$  should be expected. As expected, the observed  $U_{0.10}$  did not differ between the DHF1 and DHF2 lines for yield and plant height in either population. Despite the higher  $V_G$  for moisture among the DHF2 lines, the observed  $U_{0.10}$  for moisture did not differ significantly between the DHF1 and DHF2 lines. The difference in the observed  $U_{0.10}$  for moisture between the DHF1 and DHF2 lines was  $222.4 - 223.1 = -0.7 \text{ g kg}^{-1}$  in Population A, and  $195.4 - 197.8 = -2.4 \text{ g kg}^{-1}$  in Population B (Table 1). For comparison, the difference in the predicted  $U_{0.10}$  for moisture between the DHF1 and

DHF2 lines was  $-3.0 \text{ g kg}^{-1}$  in Population A and  $-4.4 \text{ g kg}^{-1}$  in Population B (Table 1). The differences for moisture in both observed  $U_{0.10}$  and predicted  $U_{0.10}$  were therefore very small (i.e., less than one half a percentage point of moisture at harvest).

In this study, any significant differences in  $V_G$  between the DHF1 and DHF2 lines therefore did not translate to a significant difference in  $U_{0.10}$ . This result was consistent with previous findings in maize in which additional recombination via random mating in four populations produced no substantial increases in genetic gain (Covarrubias-Prieto, 1987), and random mating in a  $BC_1$  generation did not lead to a higher  $V_G$  or  $U_{0.10}$  (Arbelbide and Bernardo, 2004).

In conclusion, the results from this study indicated that the decision of whether to induce doubled haploids from  $F_1$  or  $F_2$  plants in maize needs to be made from considerations other than the mean,  $V_G$ , or  $U_{0.10}$  in the resulting lines. These considerations would include potentially larger genetic gains in the long term with DHF2 lines than with DHF1 lines (Bernardo, 2009), the ability to select among individual  $F_2$  plants prior to inducing DHF2 lines, or the shorter time needed to produce DHF<sub>1</sub> lines.



Table 1: Trait means, genetic variance ( $V_G$ ), residual variance ( $V_R$ ), heritability ( $h^2$ ), and mean of best 10% of testcrosses ( $U_{0.10}$ ) in DHF1 and DHF2 maize populations

Trait	Generation	Population A				
		Mean	$V_G$	$V_R$	$h^2$	$U_{0.10}$
Yield	DHF1	12.72	0.42 (0.31, 0.56) <sup>b</sup>	2.12	0.54 <sup>a</sup>	13.50
Yield	DHF2	12.78	0.37 (0.25, 0.51)	2.25	0.49	13.57
Moisture	DHF1	232.9	51.5 (36.7, 68.0)	226.1	0.58	223.1
Moisture	DHF2	233.1	81.4 (60.5, 103.6)	195.5	0.71	222.4
Plant height	DHF1	238.3	103.1 (82.3, 124.6)	151.3	0.80	252.9
Plant height	DHF2	237.8	91.9 (72.7, 114.6)	131.4	0.81	252.6
		Population B				
		Mean	Mean	Mean	Mean	Mean
Yield	DHF1	13.83	13.83	13.83	13.83	13.83
Yield	DHF2	13.96	13.96	13.96	13.96	13.96
Moisture	DHF1	214.4	214.4	214.4	214.4	214.4
Moisture	DHF2	213.1	213.1	213.1	213.1	213.1
Plant height	DHF1	264.0	264.0	264.0	264.0	264.0
Plant height	DHF2	263.4	263.4	263.4	263.4	263.4

<sup>a</sup> All estimates of  $h^2$  were significant ( $P = 0.05$ )

<sup>b</sup> Lower and upper limits of 95% confidence intervals in parenthesis

Table 2: Mean and variance of the number of recombinations per chromosome and the frequency of inheriting an entire parental chromosome among maize DHF1 and DHF2 lines

Generation	Statistic	Chromosome										Total
		1	2	3	4	5	6	7	8	9	10	
DHF1	Mean	2.20	0.64	1.59	1.29	2.56	1.16	1.56	1.09	1.38	1.02	14.48
DHF2	Mean	3.28	0.89	2.40	2.17	3.53	1.76	2.21	1.75	2.04	1.35	21.38
DHF1	Variance	1.96	0.52	1.17	0.88	2.65	0.84	1.53	0.84	1.45	0.75	14.02
DHF2	Variance	3.34	0.75	2.26	1.82	3.72	1.53	2.51	1.50	2.51	0.98	29.06
DHF1	Frequency of no recombinations	0.08	0.49	0.16	0.21	0.09	0.25	0.19	0.28	0.23	0.31	0.23
DHF2	Frequency of no recombinations	0.05	0.38	0.08	0.09	0.04	0.18	0.10	0.16	0.12	0.17	0.14

Table 3: Repulsion:coupling ratios (95% confidence intervals in parenthesis) among DHF1 and DHF2 maize lines

Trait	Generation	Repulsion:coupling (R:C) ratio	
		Population A	Population B
Yield	DHF1	0.70 (0.56, 0.82)	0.52 (0.42, 0.63)
Yield	DHF2	0.60 (0.49, 0.71)	0.52 (0.43, 0.65)
Moisture	DHF1	0.44 <sup>a</sup> (0.33, 0.59)	0.59 <sup>a</sup> (0.50, 0.69)
Moisture	DHF2	0.31 (0.26, 0.36)	0.35 (0.29, 0.42)
Plant height	DHF1	0.54 (0.42, 0.67)	0.35 (0.28, 0.44)
Plant height	DHF2	0.49 (0.40, 0.58)	0.29 (0.22, 0.37)

<sup>a</sup>: Differences between the means of the DHF1 and DHF2 lines were significant ( $P = 0.05$ ) only for moisture in each population

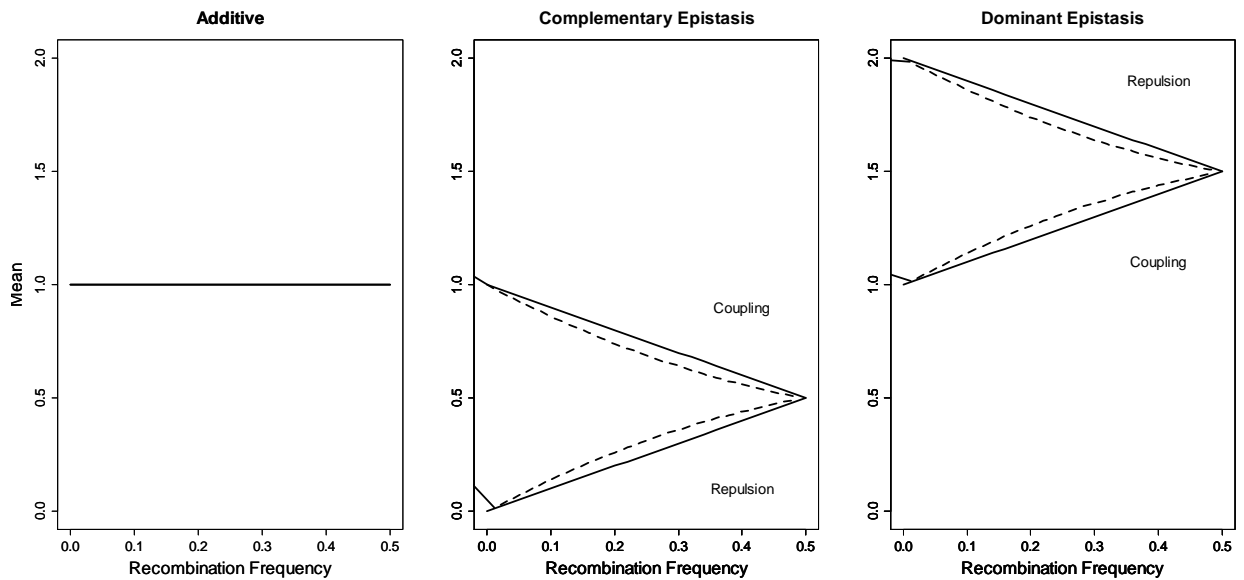
Table 4: Estimated repulsion:coupling (R:C) ratios and correlation between underlying and estimated R:C ratios among simulated DHF1 and DHF2 lines

Underlying R:C ratio	DHF1 lines						DHF2 lines					
	$h^2 = 0.3^a$			$h^2 = 0.8$			$h^2 = 0.3$			$h^2 = 0.8$		
	10 QTL	19 QTL	28 QTL	10 QTL	19 QTL	28 QTL	10 QTL	19 QTL	28 QTL	10 QTL	19 QTL	28 QTL
0:9	0.095	0.093	0.093	0.173	0.158	0.144	0.114	0.112	0.113	0.197	0.168	0.171
2:7	0.105	0.091	0.087	0.186	0.160	0.150	0.125	0.102	0.107	0.215	0.175	0.172
4:5	0.102	0.092	0.104	0.197	0.165	0.157	0.120	0.118	0.108	0.242	0.198	0.193
6:3	0.098	0.092	0.100	0.227	0.180	0.170	0.126	0.118	0.119	0.248	0.215	0.209
8:1	0.111	0.109	0.101	0.254	0.207	0.200	0.129	0.125	0.123	0.283	0.248	0.236
Correlation	0.61	0.68	0.67	0.98	0.92	0.94	0.83	0.79	0.72	0.98	0.98	0.97

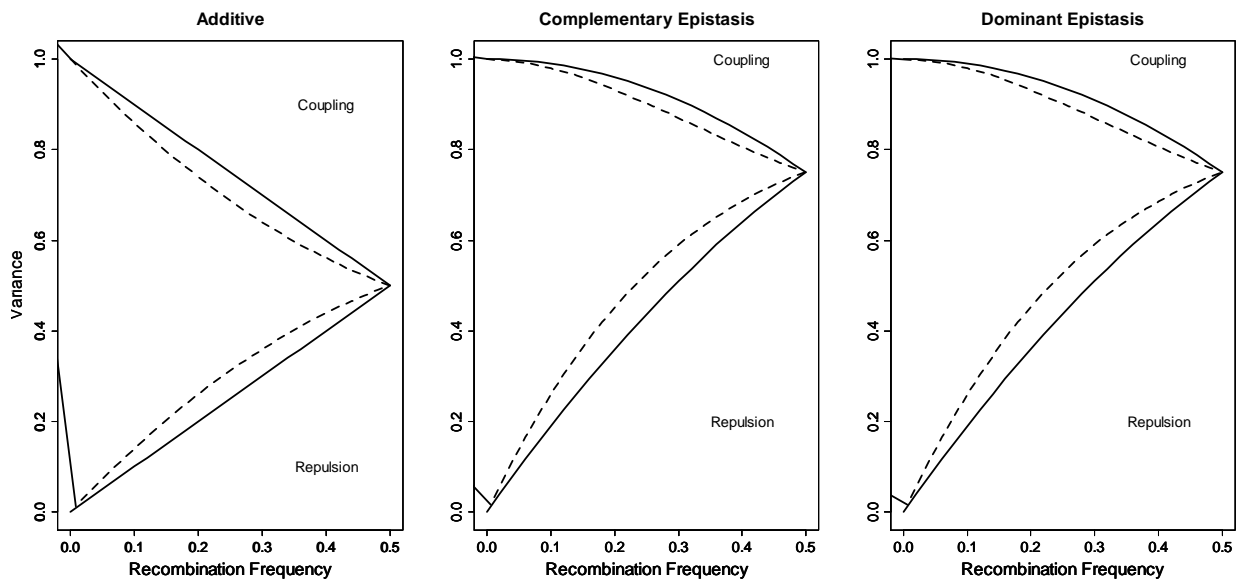
<sup>a</sup>: Results (not shown) for a heritability of  $h^2 = 0.50$  were intermediate to the results for  $h^2 = 0.30$  and  $0.80$

Figure 1: Two-locus models for the expected mean and genetic variance among DHF1 lines (solid line) and DHF2 lines (dashed line). The genotypic values were 2 for the *A-B-* genotype and 0 for all other genotypes with complementary epistasis; and 0 for the *aabb* genotype and 2 for all other genotypes with dominant epistasis

## Mean



## Genetic variance



## **Chapter 2: Genomewide Selection with Biallelic versus Triallelic Models in Three-way Maize Populations**

While single nucleotide polymorphism (SNP) markers are typically biallelic, quantitative trait loci (QTL) may have three alleles per locus in three-way populations. Our objective in this study was to determine if multi-allelic markers or haplotypes improve the prediction accuracy of genomewide selection in three-way breeding populations. Simulated and empirical maize (*Zea mays* L.) double haploid populations were used to compare a Biallelic model, Marker Interval model (which used adjacent markers to create haplotypes), and Allele Phasing model (which inferred triallelic markers from parental SNP data). The simulation experiments differed in the number of QTL (10, 40, or 100), heritability (0.30, 0.50, or 0.80), and sizes of allelic effects. Four empirical three-way populations were phenotyped at 4–7 locations between 2012 and 2015 and were genotyped with 356–960 polymorphic SNP markers. Genomewide marker effects were obtained by ridge regression-best linear unbiased prediction. In the simulation experiments, differences in prediction accuracy were less than 0.01 among the Biallelic, Marker Interval, and Allele Phasing models. For grain yield, moisture, and test weight in the four maize populations, the differences in predictive ability among the three models were nonsignificant ( $P = 0.05$ ). Further simulations showed that the small or nonsignificant differences in prediction accuracy were caused by large linkage blocks found among inbreds, particularly double haploids. Overall, we recommend the Marker Interval model in three-way populations because of its simplicity, similar prediction accuracy, and theoretical advantage over the two other models.

## Introduction

Maize (*Zea mays* L.) breeders usually create breeding populations by crossing two inbreds (A and B). Within an A/B biparental cross, genomewide selection (or genomic selection) can be routinely performed via the general combining ability (GCA) model (Jacobson et al., 2014). The GCA model involves pooling A/\* and B/\* half-sib populations, where \* is an inbred from the same heterotic group as A and B, into a training population to predict the performance of progeny within the A/B population. Separate prediction equations are therefore developed for each A/B population. Results for 969 maize biparental populations showed that, on average, genomewide selection via the GCA model led to about 85% of the gains achieved by phenotypic selection at about 25% of the cost (Jacobson et al., 2014).

Additionally, maize breeders occasionally create breeding populations from three or more inbreds. A compilation of maize Plant Variety Protection certificates between 1980 and 2004 reported that elite inbreds were derived from three-way populations 5% of the time (Mikel and Dudley, 2006). A maize three-way population is developed by crossing an A/B F<sub>1</sub> with a third inbred (C) to create an (A/B)//C population. Inbreds are then developed either through self-pollination or through induction of doubled haploids (DH). A maize three-way population is typically created with two breeding goals in mind: to improve a debilitating trait of one parent, or to introduce new genetic variation while maintaining local adaptation. For the first goal, if A is high-yielding but is extremely susceptible to a disease, a maize breeder may cross A to an inbred (B) that is closely related to A and is resistant to the disease. The A/B F<sub>1</sub> can be crossed to a third inbred (C) that is elite and not closely related to either A or B, and progeny can then be developed from the (A/B)//C

cross. For the second goal, a maize breeder may wish to widen the genetic base by incorporating a late-maturing, less-adapted inbred (D) into the breeding germplasm. Crossing inbreds C and D would most likely result in progeny with too late a maturity. The breeder could then decide to use the earliest maturing line in the current germplasm as the third parent (E) to create the (D/E)//C cross.

A three-way population introduces complexities in genomewide selection. Single nucleotide polymorphism (SNP) markers used in genomewide selection are typically biallelic. In a biparental population with inbred parents, the biallelic SNP markers are biologically reflective because there can only be two alleles at each of the underlying quantitative trait loci (QTL). But in a three-way population, the biallelic SNP markers may fail to capture the variation across the three possible alleles at each QTL. While most genomewide selection models have assumed biallelic marker effects, multi-allelic models have been proposed by Calus et al. (2008), Cuyabano et al. (2014), and Da (2015). These three studies sought to incorporate haplotypes in models to predict performance in animal populations. Compared with a biallelic model, haplotype-based models led to a 3.1% improvement in prediction accuracy for milk protein in dairy cattle (*Bos taurus*) (Cuyabano et al., 2014).

To our knowledge, there have been no published studies exploring multi-allelic or haplotype genomewide selection systems in plants. Therefore, our objective in this study was to determine if multi-allelic markers improve the prediction accuracy of genomewide selection in simulated as well as empirical three-way maize populations.



## Materials & Methods

### *Biallelic Model*

In the Biallelic model, the GCA model was expanded to accommodate three-way populations while still having biallelic SNP loci. In the GCA model, effect of the SNP allele from inbred A (in the A/B population) is calculated as the mean effect of that same allele across different A/\* populations (Jacobson et al., 2014). The effect of the SNP allele from inbred B is calculated as the mean effect of that same allele across different B/\* populations. In the Biallelic model for a three-way population, the effects of the alleles from parents A, B, and C were likewise calculated as the mean effects of those alleles in the A/\*, B/\*, and C/\* populations, respectively. Marker effects were calculated for the  $N_M$  SNP markers that were polymorphic in the (A/B)//C test population, and the mean marker effects across the A/\*, B/\*, or C/\* populations were used to predict the performance of the (A/B)//C test population. To illustrate, suppose the marker genotypes at a given SNP locus were *MM* in A, *mm* in B, and *MM* in C. Further suppose there were 10 A/\* populations and 5 C/\* populations. In this situation, the effect of the *M* allele was calculated within each of the 10 A/\* crosses and 5 C/\* crosses, and the effect of *M* was calculated as the mean of the 15 within-cross effects. This mean effect was used in subsequently predicting the performance of progeny in the (A/B)//C cross.

Specifically, the performance of each of the  $N$  individuals in the (A/B)//C test population was predicted as  $\mathbf{y} = \mu + \mathbf{M}\mathbf{g}$ , where  $\mathbf{y}$  was an  $N \times 1$  vector of the predicted performance for the trait;  $\mu$  was the population mean;  $\mathbf{M}$  was an  $N \times N_M$  incidence matrix with 1 and  $-1$  for the contrasting homozygous SNP genotypes and 0 for the heterozygous

genotype or missing data; and  $\mathbf{g}$  was an  $N_M$  vector of the mean marker effects across the A/\*, B/\*, and C/\* crosses.

#### *Marker Interval Model*

The Marker Interval model, which was the “SNP2” model of Calus et al. (2008), used the interval between two adjacent markers as a unique haplotype. For any two adjacent SNPs, there were four potential genotypes in an inbred:  $M_1M_1M_2M_2$ ,  $M_1M_1m_2m_2$ ,  $m_1m_1M_2M_2$ , and  $m_1m_1m_2m_2$ . These four haplotypes were then defined as four unique alleles. Using the progeny in Fig. 2b as an example, the first and second SNP markers were used to create four alleles, the second and third SNP markers were used to create four alleles, and so on for the other SNP markers. The coding of alleles was consistent between the test and training populations. While the model can accommodate four alleles per locus, only three alleles per locus were present in each three-way test population.

In the Marker Interval model, each allele was interpreted as a presence/absence variable (where 1 indicated the presence of a given allele and 0 indicated the absence of the allele). In contrast, the Biallelic model involved a linear contrast among the three genotypes (coded as 1, 0, and  $-1$ ) at a marker locus. The Marker Interval model therefore required four effects to be estimated (one of the effects being null) per SNP locus whereas the Biallelic model required only one effect to be estimated. Allele effects were estimated in each A/\*, B/\*, and C/\* population and were averaged in the same way as in the Biallelic model. The performance of each individual in the A/B test population was predicted as  $\mathbf{y} = \mu + \mathbf{M}\mathbf{I}\mathbf{g}$ , where  $\mathbf{M}\mathbf{I}$  was an  $N \times 4N_M$  incidence matrix and  $\mathbf{g}$  was a  $4N_M \times 1$  vector of marker allele effects.

#### *Allele Phasing Model*

The Allele Phasing model inferred multiple alleles at the SNP loci themselves. The three alleles in a three-way DH population were inferred from segregating markers and the parental genotypes. Consider the hypothetical genotypes on a single chromosome in parents A, B, and C (Fig. 2a). The orange-colored SNP alleles in Fig. 2a are those that can be uniquely traced from one of the three parents. For a biallelic marker, there will always be a 1:2 ratio of parental genotypes if a marker locus is polymorphic among the three parents. Suppose parents A and B have the *MM* genotype (coded as 1) at a given SNP while parent C has the *mm* genotype (coded as -1) (Fig. 2a). The descent of the *m* allele can be traced to parent C for all DH progeny that have the *mm* genotype at this locus. But at this juncture (i.e., prior to allele phasing), the descent of the *M* allele cannot be traced to either A or B. The same logic applies to genotypes unique to parents A and parents B.

To phase the alleles onto the progeny, markers polymorphic among the three inbred parents were used and alleles that could be traced uniquely to a parent were first identified among the progeny (Fig. 2b). These alleles were tagged as *A*, *B*, and *C* (Fig. 2c). The missing alleles in the third pane of Fig. 2c were then projected (Fig. 2d) assuming no recombination. This assumption was supported by the minimal recombination in DH lines (Smith et al., 2008; Sleper and Bernardo, 2016) and the dense linkage map (mean marker spacings of 1.69 cM in the simulation experiments and 2.2 cM in the empirical experiments described later). However, marker alleles were left as missing when a recombination was identified but the parental alleles could not be identified. For example, DH line 4 of Fig. 2d has missing data for one SNP locus.

Once projected, each marker in the progeny then had the three alleles *A*, *B*, or *C*. The *A*/\* populations were used to estimate the effect of allele *A*, the *B*/\* populations were

used to estimate the effect of allele  $B$ , and the  $C/*$  populations were used to estimate the effect of allele  $C$ . The performance of the progeny in the test population was predicted as  $\mathbf{y} = \mu + \mathbf{M_P g_P}$ , where  $\mathbf{M_P}$  was an  $N \times 3N_M$  incidence matrix, and  $\mathbf{g_P}$  was a  $3N_M \times 1$  vector of allelic effects.

### *Simulation Experiments*

Simulations were conducted to assess the efficacy of the Biallelic, Marker Interval, and Allele Phasing models when three different alleles were known to segregate at each QTL. Given this objective,  $A/*$ ,  $B/*$ , and  $C/*$  crosses were not simulated for the sake of simplicity. Instead, an  $(A/B)//C$  population of size  $N$  was simulated. As described below, the performance of each DH line was then predicted from information on the remaining  $N - 1$  lines.

Each simulation experiment constituted a different combination of the number of QTL, intensity of allele effects, and heritability ( $h^2$ ). The QTL locations differed in each simulated experiment, and each experiment was repeated 100 times. In each repeat, three unique inbred parents were created resulting in a unique  $(A/B)//C$  population.

Parents  $A$ ,  $B$ , and  $C$  were polymorphic at 1000 SNP markers. The markers were evenly spaced, and the chromosome sizes corresponded to those in a maize linkage map (Senior et al., 1996). The QTL were randomly placed across the genome according to a uniform distribution. The QTL effects were additive and they varied according to a geometric series (Lande and Thompson, 1990). Furthermore, two allelic series and two allele intensities ( $a = 1$  and  $a = 4$ ) were considered. In the first allelic series, the effects of the three homozygotes were  $-a$ ,  $0$ , and  $a$ . In the second allelic series, the effects of the three

homozygotes were  $-a$ ,  $0.5a$  and  $a$ . Dominance was inconsequential because only homozygous lines were simulated.

The different alleles were alternately assigned to each of the inbred parents. For the first QTL, parent A had the favorable allele, parent B had the unfavorable allele, and parent C had the middle allele (effect of 0 or  $0.5a$ ). For the next QTL, parent B had the favorable allele, parent C had the unfavorable allele, and parent A had the middle allele. For the following QTL, parent C had the favorable allele, parent A had the unfavorable allele, and parent B had the middle allele. This pattern continued for the remaining QTL. Three different numbers of QTL were simulated (10, 40 and 100).

A total of 100 plants were simulated from the (A/B)//C cross, and a single DH line was simulated from each plant. Genotypic values were calculated for each DH line, and phenotypic values were calculated by adding random nongenetic effects to the known genotypic values. The nongenetic effects had a normal distribution with a mean of zero and a variance scaled according to  $h^2$  (0.30, 0.50, 0.80).

Genomewide marker effects were estimated by ridge regression-best linear unbiased prediction (RR-BLUP) with the R package rrBLUP (Endelman, 2011). The performance of the DH lines in the (A/B)//C population was predicted using the Biallelic, Marker Interval, and Allele Phasing models. The correlation between the predicted values and the true genotypic values ( $r_{MG}$ ) was calculated through delete-one cross validation. A 95% confidence interval was constructed for the pairwise difference in  $r_{MG}$  for the three models by resampling on the simulation repeats. Resampling with replacement was done 1000 times and, for each sample, the mean pairwise difference in  $r_{MG}$  between models was calculated. The 1000 values of the difference in  $r_{MG}$  were sorted in ascending order. The

25th sorted value corresponded to the lower limit and the 975th sorted value corresponded to the upper limit of a 95% confidence interval. The accuracy of the Allele Phasing method was tested by calculating total percentage of alleles correctly projected from the parents to the progeny. Linkage disequilibrium was calculated as the mean  $r^2$  value between adjacent SNP markers.

### *Maize Empirical Data*

Phenotypic and genotypic data were obtained for 109 DH maize populations proprietary to Syngenta Seeds, LLC. Of these populations, four three-way populations were chosen as test populations. One of the populations was derived from inbreds from the Iowa Stiff Stalk Synthetic (BSSS) heterotic group and the other three populations were derived from inbreds from the non-BSSS heterotic group. The 105 remaining populations had one parent in common with the four three-way populations and were used as the A/\*, B/\*, or C/\* populations. Each population was testcrossed to an inbred from the other heterotic group, with the same tester being used for a test population and corresponding training populations. All phenotypic data used in this study were testcross data.

Testcrosses of the populations were grown in 4–7 locations in the U.S. and southern Europe between 2012 and 2015. Data were collected for grain yield ( $\text{Mg ha}^{-1}$  at  $155 \text{ g H}_2\text{O kg}^{-1}$ ), grain moisture ( $\text{g kg}^{-1}$ ), and test weight ( $\text{kg hL}^{-1}$ ). All trials had one replication per location, and each phenotypic data point was the performance of an individual within each location. No adjustments were made for field spatial variability within each location.

Across-location least squares means were estimated for each individual in each population using the model  $y_{ij} = \mu + g_i + l_j + e_{ij}$ , where  $y_{ij}$  was the phenotypic value of the  $i$ th individual at the  $j$ th location;  $\mu$  was the grand mean;  $g_i$  was the effect of the  $i$ th individual;  $l_j$  was the effect of the  $j$ th location; and  $e_{ij}$  was the residual effect. Testcross

genetic variance ( $V_G$ ) and nongenetic variance ( $V_R$ ) were calculated from an across-locations mixed-model analysis using the lmer function in the lme4 package in R (Bates et al., 2013). The genotype-by-environment interaction variance and within-location error variance could not be estimated separately and were confounded in  $V_R$ . A likelihood ratio test was used to determine the significance ( $P = 0.05$ ) of the  $V_G$  in each population. The entry-mean  $h^2$  was estimated for each population as  $V_G/(V_G + V_R/e)$ , where  $e$  was the number of environments.

Each DH line was genotyped with 3,072 SNP markers (proprietary to Syngenta Seeds, LLC) using the Illumina GoldenGate platform. The parental inbreds and testers were also genotyped using the same 3,072 SNP markers. The numbers of polymorphic SNP loci were 356 in (P1/P2)//P3, 857 in (P4/P5)//P6, 902 in (P4/P7)//P6, and 960 in (P8/P9)//(P8/P10).

Genetic similarity was calculated between parent A and parent B ( $S_{A|B}$ ). The genetic similarity between parent C and A/B ( $S_{AC|BC}$ ) was calculated by taking the mean value of the genetic similarity between parent A and parent C and the genetic similarity of parent B and parent C. The genetic similarity was calculated as the simple matching coefficient (Sokal and Michener, 1958) across the 3,072 SNP markers.

Marker effects were calculated using the same methodology used in the simulation experiments. The predictive ability ( $r_{MP}$ ) was calculated as a correlation between marker-predicted values and phenotypic values. The test statistic  $T = [r_{MP}(N - 1)^{1/2}]/(1 - r_{MP}^2)$  was calculated for  $r_{MP}$  (Bobko, 2001). Significance tests ( $P = 0.05$ ) for the test statistic  $T$  were done using a t-test. Significance tests for the differences in  $r_{MP}$  between models were done via a standard Fisher z-transformation for correlation coefficients.

## Results and Discussion

### *Simulated Populations*

The simulation results showed very small differences in  $r_{MG}$  values among the Biallelic, Marker Interval, and Allele Phasing models (Table 5). For each of the 18 genetic models, the mean difference in  $r_{MG}$  among the Biallelic, Marker Interval, and Allele Phasing models was less than 0.01. These differences in  $r_{MG}$  were often statistically significant ( $P = 0.05$ ), with the Marker Interval model consistently ranking better (by less than 0.01) than the Biallelic model and the Allele Phasing model in terms of  $r_{MG}$ . These differences in  $r_{MG}$  were too small to be meaningful in a plant breeding program. Among the simulation parameters,  $h^2$  had the largest effect on  $r_{MG}$  values and, as expected (Daetwyler et al., 2008; Lian et al., 2014),  $r_{MG}$  values increased as  $h^2$  increased (Table 5). There was no clear difference in  $r_{MG}$  values due to allelic series, the allele-effect intensities, or number of QTL.

Across all simulated populations, the Allele Phasing algorithm correctly identified 96% of the underlying marker alleles across the 1000 loci in each three-way DH population. The remaining 4% of alleles were considered as missing and were located between recombination sites. In particular, 90% of missing data points were within 5 cM of a recombination site. These results indicated a high accuracy of the algorithm for inferring SNP alleles from parental haplotypes in DH populations.

### *Empirical Populations*

The number of A/\*, B/\*, and C/\* populations that comprised the training population for each (A/B)//C population ranged from 2 to 56 (Table 6). The size of the training



population ranged from 644 to 5,757 inbreds. Heritability in each three-way test population was significant ( $P = 0.05$ ) for each trait. The mean  $h^2$  across the four test populations was 0.51 for grain yield, 0.65 for grain moisture, and 0.61 for test weight. The mean linkage disequilibrium between adjacent markers was  $r^2 = 0.73$  across all four populations.

For each trait, the Biallelic, Marker Interval, and Allele Phasing models did not differ significantly ( $P = 0.05$ ) in  $r_{MP}$  (Table 7). The observed differences in  $r_{MP}$  were less than 0.05 in most of the population-trait combinations. For the first three populations, all  $r_{MP}$  values were significantly different from zero ( $P = 0.05$ ); however,  $r_{MP}$  was significant only for moisture in the fourth population [(P8/P9)/(P8/P10)] for the Marker Interval and Allele Phasing models (Table 7). Overall, the results from the four empirical maize populations were consistent with the results from the simulation experiments.

#### *Equivalency between the Models in Practice*

In a previous simulation study, the prediction accuracy was higher with the Marker Interval model than with the Biallelic model especially under dense genotyping (0.5–1.0 cM marker spacing) (Calus et al., 2008). In theory, the Allele Phasing model and the Marker Interval model are expected to outperform the Biallelic model in a three-way population with multiple alleles. To illustrate, assume that a QTL with multiple alleles is completely linked to a biallelic SNP marker. Parent A has the SNP genotype coded as 1 and an effect of 1; Parent B has the SNP genotype coded as -1 and an effect of 0; and Parent C has the SNP genotype coded as 1 and an effect of -1. If the biallelic SNP locus is used to estimate the effect at the underlying QTL, the estimates will be confounded with the linkage phases because Parent A and Parent C share the same SNP marker, but have contrasting alleles. In this example, the  $r_{MG}$  with the Biallelic model is only 0.17.

In contrast to these expectations, the simulation experiments and empirical maize populations showed no differences in prediction accuracy among the Biallelic, Marker Interval, and Allele Phasing models (Table 7). While the example in the previous paragraph demonstrated a low  $r_{MG}$  with the Biallelic model, a different linkage phase can lead to a high  $r_{MG}$  with the Biallelic model. To expand on the example in the previous paragraph, suppose Parent A has the SNP genotype coded as 1 and an effect (due to a perfectly linked QTL) of 1; Parent B has the SNP genotype coded as 1 and an effect of 0; and Parent C has the SNP genotype coded as -1 and an effect of -1. In this example, the  $r_{MG}$  with the Biallelic model is 0.90.

The two examples shown above might not reflect how the effects of multiple linked QTL are captured by multiple linked SNP loci. We therefore further examined the effect of linkage disequilibrium on prediction accuracy in three-way populations. In particular, we simulated 100 three-way populations for three of the genetic models (100 QTL,  $h^2 = 0.30$ , first allelic series,  $a = 1$ ; 40 QTL,  $h^2 = 0.50$ , first allelic series,  $a = 1$ ; 10 QTL,  $h^2 = 0.80$ , first allelic series,  $a = 1$ ). Instead of immediately developing DH lines from the  $F_2$  generation, the simulated  $F_2$  plants were random mated for six generations before developing DH lines. In all three genetic models, random mating drastically reduced the effectiveness of the Biallelic model but not of the Marker Interval and Allele Phasing models (Table 8). To illustrate, for the second genetic model (which involved an intermediate number of QTL and intermediate  $h^2$ ), random mating decreased the mean  $r_{MG}$  from 0.48 (Table 5) to -0.11 for the Biallelic model; from 0.49 to 0.35 for the Marker Interval model; and from 0.48 to 0.32 for the Allele Phasing model. The results suggested that high linkage disequilibrium caused the biallelic and triallelic models to be equivalent

in both the simulation studies and empirical populations, but a low linkage disequilibrium is expected to drastically reduce the effectiveness of the Biallelic model in populations with more than two alleles per QTL.

The small differences in  $r_{MG}$  between the Allele Phasing and Marker Interval models were most likely explained by the imperfect accuracy of the projection algorithm in the Marker Interval model. On average, the algorithm left 4% of the markers as having missing data. This observed level of inaccuracy was low, but it would likely be higher in populations with lower linkage disequilibrium than in DH populations. The Marker Interval model should therefore be preferred over the Allele Phasing model in three-way populations. Additionally, the Allele Phasing model requires genotypic data of the parents used in the population, but parental data might not always be available.

#### *Application*

The results indicated that maize breeders should not be deterred from genomewide selection in three-way populations. We recommend the Marker Interval model in three-way populations because of its simplicity, practical equivalency (in terms of  $r_{MG}$  or  $r_{MP}$ ) with the two other models, and theoretical advantage over the two other models. However, the number of marker effects to be calculated is larger with the Marker Interval model than with the Biallelic model. This difference may be a hindrance to the Marker Interval model when the number of SNP markers is large.

Studies will be needed to determine the efficacy of multiple-allele models in more complex populations such as a four-parent population or a synthetic population. In a rice (*Oryza sativa* L.) synthetic population, prediction accuracies with the Biallelic model ranged from 0.30 for flowering date to 0.54 for plant height (Grenier et al., 2015), and it

remains to be seen whether the Marker Interval model would increase the prediction accuracy in a synthetic. Multiple-allele models may also be needed for genomewide selection in broadbase maize populations, such as those undergoing long-term recurrent selection (Hallauer and Carena, 2012).

Table 5: Estimated prediction accuracy in simulated three-way populations.

Allelic series	Model	Heritability = 0.30					
		10 QTL		40 QTL		100 QTL	
		$a = 1$	$a = 4$	$a = 1$	$a = 4$	$a = 1$	$a = 4$
$-a, 0, a^\dagger$	Biallelic	0.34 <sup>‡</sup>	0.32	0.36	0.32	0.32	0.35
$-a, 0, a$	Marker Interval	0.33	0.32	0.36	0.33	0.32	0.36
$-a, 0, a$	Allele Phasing	0.33	0.32	0.36	0.33	0.32	0.35
$-a, 0.5a, a$	Biallelic	0.33	0.33	0.31	0.31	0.31	0.31
$-a, 0.5a, a$	Marker Interval	0.34	0.34	0.33	0.31	0.31	0.31
$-a, 0.5a, a$	Allele Phasing	0.33	0.33	0.32	0.31	0.31	0.31
		Heritability = 0.50					
		10 QTL		10 QTL		10 QTL	
		$a = 1$	$a = 1$	$a = 1$	$a = 1$	$a = 1$	$a = 1$
$-a, 0, a^\dagger$	Biallelic	0.45	0.45	0.45	0.45	0.45	0.45
$-a, 0, a$	Marker Interval	0.45	0.45	0.45	0.45	0.45	0.45
$-a, 0, a$	Allele Phasing	0.45	0.45	0.45	0.45	0.45	0.45
$-a, 0.5a, a$	Biallelic	0.46	0.46	0.46	0.46	0.46	0.46
$-a, 0.5a, a$	Marker Interval	0.47	0.47	0.47	0.47	0.47	0.47
$-a, 0.5a, a$	Allele Phasing	0.47	0.47	0.47	0.47	0.47	0.47
		Heritability = 0.80					
		10 QTL		10 QTL		10 QTL	
		$a = 1$	$a = 1$	$a = 1$	$a = 1$	$a = 1$	$a = 1$
$-a, 0, a^\dagger$	Biallelic	0.62	0.62	0.62	0.62	0.62	0.62
$-a, 0, a$	Marker Interval	0.63	0.63	0.63	0.63	0.63	0.63
$-a, 0, a$	Allele Phasing	0.62	0.62	0.62	0.62	0.62	0.62
$-a, 0.5a, a$	Biallelic	0.61	0.61	0.61	0.61	0.61	0.61
$-a, 0.5a, a$	Marker Interval	0.61	0.61	0.61	0.61	0.61	0.61
$-a, 0.5a, a$	Allele Phasing	0.60	0.60	0.60	0.60	0.60	0.60

<sup>†</sup> Coded genotypic value of the favorable homozygote

<sup>‡</sup> The least significant difference ( $P = 0.05$ ) across all simulation experiments was less than 0.01.

Table 6: Summary statistics for assessing the predictive ability of genomewide selection in each of four three-way maize test populations.

Three-way test population									Training population						
Cross	Tester	$S_{A/B}^{\dagger}$	$S_{AC BC}^{\ddagger}$	$N^{\S}$	Locations	Heritability			A/* <sup>¶</sup>	B/*	C/*	$N_{Train}^{\#}$	Heritability <sup>††</sup>		
						Yield	Moisture	Test weight					Yield	Moisture	Test weight
(P1/P2)//P3	T1	0.81	0.88	45	5	0.596	0.79	0.779	2	7	56	5757	0.40 (0.13, 0.78) <sup>‡‡</sup>	0.67 (0.37, 0.93)	0.73 (0.17, 0.91)
(P4/P5)//P6	T2	0.73	0.75	98	5	0.491	0.689	0.598	5	7	3	962	0.49 (0.23, 0.68)	0.71 (0.29, 0.86)	0.67 (0.45, 0.90)
(P4/P7)//P6	T2	0.69	0.74	98	7	0.347	0.61	0.463	5	2	3	644	0.48 (0.23, 0.68)	0.71 (0.29, 0.86)	0.63 (0.45, 0.90)
(P8/P9)// (P8/P10)	T3	0.78	0.63	34	4	0.587	0.496	0.211	3	10	2	1694	0.48 (0.30, 0.65)	0.55 (0.24, 0.8)	0.60 (0.36, 0.81)

<sup>†</sup>  $S_{A/B}$ , genetic similarity between parent A and parent B of the three-way cross

<sup>‡</sup>  $S_{AC|BC}$ , mean of the genetic similarity between parent A and parent C and the genetic similarity between parent B and parent C

<sup>§</sup> N, number of doubled haploid lines in the three-way test population

<sup>¶</sup> A/\*, B/\*, C/\*, the number of populations where A, B, and C were the inbred parents of the (A/ B)// C cross and \* was an inbred from the same heterotic group as A, B, or C

<sup>#</sup>  $N_{Train}$ , total number of doubled haploid lines in the training population

<sup>††</sup> All heritability estimates in the test population were significantly different from zero ( $P = 0.05$ )

<sup>‡‡</sup> Median and range (in parentheses) of heritability in the A/\*, B/\*, and C/\* crosses. All heritability estimates were significantly different from zero ( $P = 0.05$ )

Table 7: Predictive ability in maize three-way test populations for yield, moisture, and test weight.

Trait	Model	Three-way test population			
		(P1/P2)//P 3	(P4/P5)//P 6	(P4/P7)//P 6	(P8/P9)//(P8/P10 )
Yield	Biallelic	0.51	0.35	0.47	0.22 <sup>NS</sup>
	Marker Interval	0.50	0.38	0.47	0.18 <sup>NS</sup>
	Allele Phasing	0.46	0.32	0.47	0.19 <sup>NS</sup>
Moisture	Biallelic	0.24	0.22	0.45	0.23 <sup>NS</sup>
	Marker Interval	0.32	0.21	0.46	0.30
	Allele Phasing	0.38	0.22	0.46	0.28
Test weight	Biallelic	0.23	0.60	0.26	0.23 <sup>NS</sup>
	Marker Interval	0.26	0.61	0.26	0.23 <sup>NS</sup>
	Allele Phasing	0.26	0.60	0.23	0.16 <sup>NS</sup>

<sup>NS</sup> Not significantly different from zero ( $P = 0.05$ ). All other estimates of  $r_{MP}$  were significant. No models were significantly ( $P = 0.05$ ) different from each other within each trait-population combination.

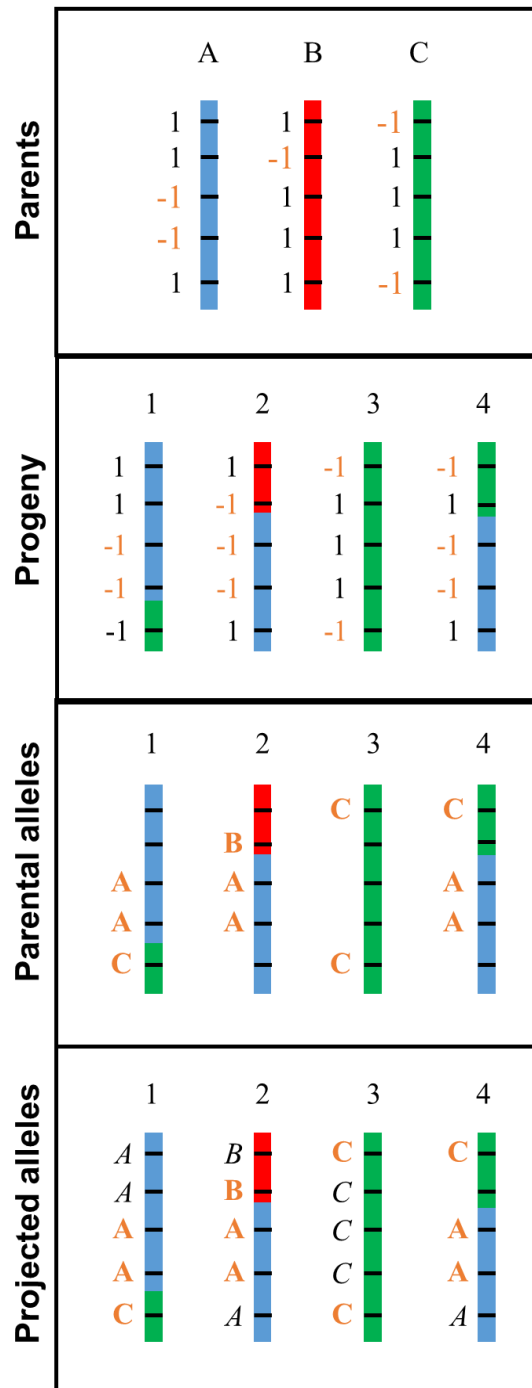
Table 8: Estimated prediction accuracy in simulated three-way population after six generations of random mating, under three genetic models that differed in heritability ( $h^2$ ) and number of quantitative trait loci (QTL).

Model	100 QTL, $h^2 = 0.30$	40 QTL, $h^2 = 0.50$	10 QTL, $h^2 = 0.80$
Biallelic	0.01	-0.11	0.00
Marker Interval	0.24	0.35	0.44
Allele Phasing	0.23	0.32	0.42

<sup>†</sup> The least significant difference ( $P = 0.05$ ) across all simulation experiments was less than 0.01



Figure 2: Allele Phasing in a three-way population: (A) inbred parents are genotyped with biallelic SNP markers; (B) inbred progeny of the three-way population are genotyped with biallelic SNP markers; (C) genotypes for a three-allele model are inferred in the three-way progeny according to allelic descent; and (D) final genotypes are projected based on flanking marker information, assuming no double recombination.



### Chapter 3: Genomewide selection for unfavorably correlated traits in maize

Correlated traits with unfavorable relationships have often hindered plant breeding efforts. Genomewide markers may help untangle unfavorable trait correlations. Our objectives were to determine (i) if genomewide markers can be used to partition trait effects into independent and correlated portions, (ii) and if selection on the independent portion will be more effective than selection on the entire trait. We used two biparental maize (*Zea mays* L.) populations to compare a standard genomewide selection model (Control Model) with a genomewide selection model that selects on the independent portion of a trait (Independent Model). We conducted genomewide selection for four cycles on these two maize populations using these two different models. In Population 1, responses to selection with the Independent Model (versus the Control Model, in parentheses) were 1.03 Mg ha<sup>-1</sup> (versus 1.40 Mg ha<sup>-1</sup>) for grain yield, 0.38 g kg<sup>-1</sup> (versus -7.98 g kg<sup>-1</sup>) for moisture, and 6.50 cm (versus 18.75 cm) for plant height. Overall, the responses were not significantly different ( $P = 0.05$ ) between the Independent Model and the Control model at each cycle of genomewide selection. The nonsignificant differences in selection responses were consistent with the low proportions ( $R^2 = 1$  to 14%) of the trait variation that was accounted for by the independent portion of the trait. We conclude that separating quantitative traits into correlated and independent portions is infeasible, mostly likely because of the complicating factors of linkage and pleiotropy.

## Introduction

Plant breeders typically select for several quantitative traits at a time. These traits are often correlated, both beneficially and adversely, amongst each other. For example, maize (*Zea mays*) biparental populations tend to have positive, unfavorable relationships between grain yield and plant height and between grain yield and moisture (Chi et al., 1969; Ross, 2002; Combs and Bernardo, 2013; Ziyomo and Bernardo, 2013). Achieving selection gains for adversely correlated traits can be difficult. One method to select for unfavorably correlated traits is to apply a selection index that accounts for both positive and negative covariances among the traits (Hazel, 1943; Smith, 1936; Baker, 1986).

Genomewide selection is a breeding strategy that leverages the use of phenotypic data along with marker data to predict the performance of non-phenotyped individuals (Meuwissen, et al., 2001). Genomewide selection studies have largely treated traits individually and have ignored their correlations (Beyene et al., 2012; Massman et al., 2013). In animal breeding, studies have focused on utilizing correlated traits in multivariate models to improve the accuracy of predictions (Aguilar et al., 2011; Calus and Veerkamp, 2011; Guo et al., 2014; Hayashi and Iwata, 2014). In plant breeding, studies have focused on using correlated, highly heritable traits to be used in indirect selection for disease resistance (Jia and Jannick, 2012; Rutkoski et al., 2012; Bao et al., 2015) or drought tolerance (Ziyomo and Bernardo, 2013). Additionally, genomewide selection for several traits in plants has involved predicting the performance for each trait individually, then combining the individual-trait predictions in a selection index (Massman et al., 2013; Combs and Bernardo, 2013; Beyene et al., 2015). In hybrid rye (*Secale cereale* L.), the accuracy was higher when a multiple-trait selection index was predicted rather than when

individual traits were first predicted and later combined into a selection index (Schulthess et al., 2015). Overall, studies in both plants and animals suggest that multiple-trait genomewide selection models will be more effective than predicting the performance of one trait at a time.

In this study, our hypothesis was that correlated traits can be modeled as having two separate components: a portion independent of another trait and a portion correlated with another trait. For example, there will be loci that affect grain yield and are physiologically independent of plant height so that they do not increase plant stature. In contrast, there will be loci that affect grain yield and that also increase plant height. These loci could exhibit by either linkage or pleiotropy. For example, the *teosinte branched1 (tb1)* gene in maize has pleiotropic effects on tillering, plant architecture (internode length, branch length, and the number of nodes) in the upper most branch, and the number of kernels per row (Clark et al., 2006). On the other hand, quantitative trait loci (QTL) associated with domestication on chromosome 5 most likely comprised multiple linked genetic factors (Lemmon and Doebley, 2014). Using nearly isogenic recombinant inbred lines, the authors were able to separate a chromosomal region associated with domestication into a QTL controlling kernel row number and a QTL controlling plant height.

In genomewide selection, the effects of genomewide molecular markers are estimated by regressing phenotypic values on molecular marker data. Molecular markers are in linkage disequilibrium with the QTL involved in the independent and correlated portions of a trait. Because of this association we should in theory be able to estimate the effects of the independent loci separate of the correlated loci. For example, in estimating

the marker effects for yield, we could separate the loci that affect yield and are independent of plant height from the loci that are associated with both yield and plant height. By using the independent portion of a trait to estimate genomewide marker effects, we might be able to maximize the selection efficiency of genomewide selection across several traits.

Such an approach therefore differs from the multivariate approaches previously described (Bao et al., 2015; Jia and Jannick, 2012; Rutkoski et al., 2012; Ziyomo and Bernardo, 2013), which involved the use of correlated traits to increase the prediction accuracy of a single trait. Our objectives in this study were to determine (i) if quantitative traits can be divided into independent and correlated portions using genomewide marker effects, (ii) and if selection using genomewide marker effects from the independent portion of quantitative traits is more effective than selection using genomewide marker effects calculated for the entire trait.

## **Materials and Methods**

### *Training Data*

Two populations from Syngenta Seeds, LLC were studied. The two populations were derived from three Lancaster SureCrop inbred parents (Population 1 was derived from the cross A x B and Population 2 was derived from the cross A x C). Each population was crossed to a different B73-type tester. Population 1 initially included 613 F<sub>1</sub> and F<sub>2</sub> induced double haploid lines, and was previously described in detail (Sleper and Bernardo, 2016). All phenotypic and genotypic analyses and results for Population 1 were reported by Sleper and Bernardo (2016). Population 1 was evaluated at six locations in Italy (Moscazzano, Torre di Mosto, and Casale Monferrato) and Spain (Palma del Río, Posadas, and Arroyo

de San Serván) in 2013. Population 2 initially included 236 F<sub>5</sub> lines and was evaluated at seven locations in the United States (Brook, Indiana; New Bedford, Illinois; Atlantic, Eagle Grove, Slater, and Stanwood, Iowa; Foster, Nebraska) in 2011. For both populations, the entries were grown in two-row plots, each row 6.1 m long and spaced 76 cm apart, at a plant population density of 89,000 plants per hectare. Data were collected for plant height (cm), yield (Mg ha<sup>-1</sup> at 155 g H<sub>2</sub>O kg<sup>-1</sup>) and moisture (g kg<sup>-1</sup>).

Each doubled haploid (DH) line from Population 1 and each F<sub>5</sub> line from Population 2 was genotyped with 3072 SNP markers (proprietary to Syngenta Seeds, LLC) on the Illumina GoldenGate platform. After removal of low quality data using the criteria given by Sleper and Bernardo (2016), 552 double haploid lines and 725 segregating SNP markers were used in the final analysis in Population 1, while 192 F<sub>5</sub> lines and 679 segregating SNP markers were used in Population 2.

For both populations, phenotypic data were analyzed within each population using a linear model including the grand mean, location effect, line effect, and residual effect with each effect assumed as random. The variance components for testcross genetic variance ( $V_G$ ) and nongenetic variance ( $V_R$ ) were calculated from across-locations analysis using the “lmer” function in the “lme4” package (Bates et al. 2013). Because the experiments used a single replication, the genotype-by-environment interaction variance and within-location error variance could not be estimated separately and were confounded in  $V_R$ . A likelihood ratio test was used to determine the significance ( $P = 0.05$ ) of the  $V_G$  in each population. The entry-mean heritability was estimated for each population as  $h^2 = V_G/(V_G + V_R/e)$ , where  $e$  was the number of environments. Within each population,

phenotypic correlations were estimated from the across-locations mean performance of the lines for each trait.

Genomewide marker effects for genomewide selection were estimated by ridge regression-best linear unbiased prediction with the R package rrBLUP (Endelman, 2011). Genotypic correlations were estimated using two methods: analysis of covariance as described in Falconer and Mackay (1996) and as the correlation between the marker effects for a pair of traits as described by Ziyomo and Bernardo (2013). Phenotypic covariances between two traits were calculated as  $\sigma_{P1,P2} = \rho_{P1,P2}\sigma_{P1}\sigma_{P2}$ , where  $\rho_{P1,P2}$  was the phenotypic correlation between two traits;  $\sigma_{P1}$  was the standard deviation of the phenotypic least-square means of the first trait; and  $\sigma_{P2}$  was the standard deviation of the phenotypic least-square means of the second trait. Genotypic covariances between two traits were calculated as  $\sigma_{G1,G2} = \rho_{G1,G2}\sigma_{G1}\sigma_{G2}$ , where  $\rho_{G1,G2}$  was the genotypic correlation between two traits;  $\sigma_{G1}$  was the square root of the genetic variance ( $V_G$ ) for the first trait; and  $\sigma_{G2}$  was the square root of the  $V_G$  of the second trait. Significance tests ( $P = 0.05$ ) for all correlations reported in this study and were done via a standard Fisher z-transformation for correlation coefficients.

#### *Genomewide Selection Models*

Two genomewide selection models were used in this study: a standard genomewide selection model (Control Model), and a model that accounted for trait correlations (Independent Model). In the Control Model, marker effects were estimated using the model  $\mathbf{y} = \mathbf{1}\mu + \mathbf{M}\mathbf{g}$ , where  $\mathbf{y}$  was a  $N \times 1$  vector of the phenotypic values of a trait;  $N$  was the number of individuals in the training population;  $\mathbf{1}$  was an  $N \times 1$  vector with elements equal to 1;  $\mu$  was the population mean;  $\mathbf{M}$  was a  $N \times N_M$  genotypic incidence matrix;  $N_M$  was the

number of polymorphic markers; and  $\mathbf{g}$  was a  $N_M \times 1$  vector of marker effects. Marker effects were estimated for each trait separately.

The purpose of the Independent Trait model was to use estimated genotypic correlations to partition phenotypic values into uncorrelated and independent portions. The independent proportion of the trait was estimated as  $\mathbf{y}_n = \mathbf{y} - \mathbf{y}_R$ , where  $\mathbf{y}_n$  was an  $N \times 1$  vector of the independent portion of the trait;  $\mathbf{y}$  was an  $N \times 1$  vector of phenotypic values for the trait; and  $\mathbf{y}_R$  was an  $N \times 1$  vector of the correlated effects of the trait. The correlated portion of the trait was estimated as  $\mathbf{y}_R = \mathbf{T}\mathbf{P}^{-1}\mathbf{g}$ , where  $\mathbf{T}$  was a  $N \times (N_T - 1)$  matrix of the phenotypic values of the other traits;  $N_T$  was the number of traits;  $\mathbf{P}$  was a  $(N_T - 1) \times (N_T - 1)$  phenotypic covariance matrix; and  $\mathbf{g}$  was a  $(N_T - 1) \times 1$  vector of genotypic covariances between the trait of interest and the other traits used in the model (Baker, 1986). Genomewide marker effects were estimated for the independent portion of each trait using the model  $\mathbf{y}_n = \boldsymbol{\mu} + \mathbf{M}\mathbf{g}$ , where  $\mathbf{g}$  was the vector of marker effects for the Independent Model. In this study, yield was adjusted for both moisture and plant height; moisture was adjusted for both yield and plant height; and plant height was adjusted for both yield and moisture. For each trait, selection was based on the adjusted (independent) portion. The variance of each trait component was calculated using  $\sigma_y^2 = \sigma_{y_n}^2 + \sigma_{y_R}^2 + 2\sigma_{y_n, y_R}$  where  $\sigma_y^2$ ,  $\sigma_{y_n}^2$ , and  $\sigma_{y_R}^2$  were calculated as the variance of each of the trait components individually. The proportion of the total variance explained by the independent and the correlated components was calculated as  $R_{\text{Independent}}^2 = \sigma_{y_n}^2 / \sigma_y^2$  and  $R_{\text{Correlated}}^2 = \sigma_{y_R}^2 / \sigma_y^2$ .

In both models, the calculated marker effects were used to predict the genotypic value of each candidate. A selection index that combined information on the three traits ( $i = 1$  to 3) was calculated as  $I = \sum b_i g_i$  where  $b_i$  was the weight for the  $i$ th trait and  $g_i$  was the



estimated genotypic value of a candidate for each trait according to the model used (Control Model or Independent Model). Each  $b_i$  value was calculated as  $w_i/\sigma_{y(i)}$ , where  $w_i$  was the weight for the trait (0.60 for yield, -0.15 for moisture, and -0.25 for plant height) and  $\sigma_y$  was the trait standard deviation.

### *Recurrent Selection*

For Population 1, the first cycle of selection (Cycle 0) consisted of 540  $F_2$  plants. For Population 2, 12 top-performing (based upon the selection index described above)  $F_5$  plants were selected and random mated using a diallel scheme. The progeny from each diallel cross were equally sampled to create a balanced bulk of 540 seeds for Cycle 0. For both populations, Cycle 0 was randomly split into two halves, each with 270 plants each. The Control Model was used for selection in the first half of the population while the Independent Model was used for selection in the second half.

Seedling DNA was extracted from leaf punches. Population 1 was genotyped with 519 of the 725 training polymorphic SNP markers and Population 2 was genotyped with 418 of the 679 training polymorphic SNP markers using TaqMan® Genotyping Assays. The SNP markers that were excluded were not adaptable to the TaqMan® Genotyping platform. For each model and population, marker effects were calculated using these sets of TaqMan® SNP markers and the original phenotypic data. Using the selection index, plants were ranked within each population and the top 12 plants were randomized and then mated by chain crossing (the first plant pollinated the second plant, the second plant pollinated the third plant, and so on). The harvested seed from these 12 plants was bulked and planted to make up Cycle 1.

The experiment was carried out for three more rounds of selection resulting in four total cycles (Cycle 0 to Cycle 4). The population size (270 plants), number of selected plants (12), SNP markers used; and SNP effects remained constant across cycles.

#### *Response to Selection*

For both Population 1 and Population 2, seed bulks were created from 100 kernels of each cycle resulting in nine seed bulks (Cycle 0, Control Model Cycles 1-4, and Independent Model Cycles 1-4). For Cycle 0, 100 F<sub>2</sub> kernels were bulked for Population 1 while 100 random mated F<sub>5</sub> kernels were bulked for Population 2. For Population 1, the bulks were grown in Grenaros, Chile during the winter season between 2015 and 2016, and were crossed to the tester used for the training population. For Population 2, Cycle 1 and Cycle 2 bulks were grown during the summer season of 2015 in Slater, Iowa, whereas Cycle 0, Cycle 3, and Cycle 4 were grown during the summer season of 2016 in Slater, Iowa. The Population 2 entries were testcrossed to the tester used for the training population.

During the summer of 2016, the nine entries from Population 1 were planted in yield trials at the same six locations used for the 2013 training population. At each location, five replications of each testcrossed seed bulk (except Cycle 0 which averaged 1.6 replications and Correlated Model Cycle 2 which average 3.6 replications due to low seed quantity) were grown with four repeated check hybrids in a randomized complete block design. Plot sizes and plant population densities were the same as those used for the training population. Plant height was measured at four locations, while yield and moisture were measured at all six locations.

For Population 2, the nine entries will be planted in yield trials in 2017 at the same locations used for the 2011 training population (as described in the Materials and Methods). The yield trials were delayed by a year because testcrosses made during the winter season of 2015-2016 were unsuccessful. The results for response to selection are still unavailable for Population 2, but will be available in the fall of 2017.

For each entry, least-square means were calculated using the R package “lsmeans” (Lenth and Hervé, 2015). The selection index was applied to the least-square means. To compare the rate of response between the Control Model and Independent Model, a linear regression model was fitted to the least-square means and the cycles. For each trait and for the selection index, regression coefficients were obtained for each model and tested for significance using a t-test ( $P = 0.05$ ). Additionally, a t-test ( $P = 0.05$ ) was used to compare the cycle means for each model. With the locations being assumed as random, the error term used in the t-test was the entry x location mean squares from analysis of variance.

To calculate changes between the cycles at the marker level, SNP markers from Cycle 1, Cycle 2, and Cycle 3 were used. The mean marker heterozygosity was calculated and used to calculate the inbreeding coefficient as  $1 - H_{Cn}/H_{C0}$ ; where  $H_{Cn}$  is the mean marker heterozygosity of the Cycle  $n$  and  $H_{C0}$  is the mean marker heterozygosity of Cycle 0. To test the difference in inbreeding between both models, bootstrapping ( $P = 0.05$ ) was conducted by sampling the chromosomes of each plant with replacement. For each cycle and population, the mean heterozygosity was calculated and the difference between the Control and the Independent models was reported.

## Results and Discussion

### *Training Population*

As previously reported (Sleper and Bernardo, 2016), estimates for  $V_G$  within Population 1 were significant for grain yield, grain moisture, and plant height. Additionally, estimates for  $V_G$  within Population 2 were significant for grain yield, grain moisture, and plant height (Table 9). Heritability was moderate to high ( $h^2 = 0.55\text{--}0.81$ ) for these three traits in both populations.

The genotypic correlations calculated from phenotypic data and from genomewide marker effects were not significantly different ( $P = 0.05$ ) from each other. The results from this study therefore support previous findings that the correlation between genomewide marker effects for separate traits can be used as an alternative way to calculate the genotypic correlation (Ziyomo and Bernardo, 2013).

For both populations, the primary trait relationship of concern was that between plant height and grain yield. The phenotypic and genotypic correlations between grain yield and plant height were strong and unfavorable in Population 1 and were moderate and unfavorable in Population 2. Additionally, the correlations between grain yield and moisture were moderate and favorable in Population 1 and were negligible in Population 2. Lastly, the correlations between plant height and grain moisture were negligible for both populations (Table 10). Overall, the correlations observed for both of these populations were consistent with previous estimates. The genetic correlation between grain yield and plant height was 0.40 in Reid Yellow Dent (Chi et al., 1969), ranged from 0.10 to 0.40 in an  $F_{2:3}$  biparental population using divergent parents derived from the Iowa Long-Ear Synthetic population (Ross, 2002), was 0.23 in the intermated B73  $\times$  Mo17 population

(Ziyomo and Bernardoi, 2013), and was 0.74-0.75 in two dwarf  $\times$  nondwarf biparental populations (Combs and Bernardo, 2013). In the intermated B73  $\times$  Mo17 population, genotypic correlations were 0.23 between grain yield and moisture and was nonsignificant ( $P = 0.05$ ) moisture and plant height. The correlations observed here and in the literature indicate that selection on grain yield without regard to plant height could hinder a maize breeding program.

The independent portion of a trait ( $y_n$ ) was highly correlated with the phenotypic values ( $y$ ) for each trait, but was not highly correlated with the other independent portions of other traits. After adjustment for trait correlations using the Independent Model, the correlation between the phenotypic value ( $y$ ) and the independent portion ( $y_n$ ) ranged from 0.93 to 1.00 for different pairs of traits in the two populations. Additionally, genomewide marker effects estimated for the Independent Model were highly correlated (0.90 to 0.99; Table 10) with marker effects estimated with the Control Model. Despite this high similarity between the models, the Independent Model genetic correlation (correlation between the marker effects of the three traits for the Independent model) was altered compared to the Control Model genetic correlation. For example, in Population 1, the Independent Model genotypic correlation between grain yield and moisture was 0.02 compared to -0.19 for the Control Model, while the genotypic correlation between grain yield and plant changed from 0.54 (Control Model) to -0.11 (Independent Model), and remained the same between plant height and grain moisture (-0.04 to -0.05) (Table 10). In Population 2, the genotypic correlations were not significant ( $P = 0.05$ ) for the Control Model (0.05) or the Independent Model (0.00) for grain yield and moisture; while the genotypic correlation changed from 0.23 (Control Model) to -0.08 (Independent Model)

for grain yield and plant height; and were not significant -0.06 (Control Model) to 0.05 (Independent Model) for plant height and grain moisture (Table 10).

The independent portions of each trait did not account for much of the overall trait variation. In Population 1,  $R^2_{\text{Independent}}$  (the proportion of the total phenotypic variance accounted for by the independent portion) was 14% for grain yield, 2% for grain moisture, and 14% for plant height; while, the  $R^2_{\text{Correlated}}$  was 69% for grain yield, 97% for grain moisture, and 71% for plant height. In Population 2,  $R^2_{\text{Independent}}$  was 4% of for grain yield, 1% for grain moisture, and 2% for plant height; while,  $R^2_{\text{Correlated}}$  was 96% for grain yield, 98% for grain moisture, and 95% for plant height. For both populations,  $R^2_{\text{Independent}} + R^2_{\text{Correlated}}$  did not equal 100% (except for grain yield in Population 2) because of the covariance between  $\mathbf{y}_n$  and  $\mathbf{y}_R$ . For Population 1,  $R^2_{\text{Independent}} + R^2_{\text{Correlated}}$  for grain yield and plant height were 83% and 85% indicating a substantial correlation between  $\mathbf{y}_n$  and  $\mathbf{y}_R$  and an underlying difficulty in separating out the two components. Overall, the high correlation between  $\mathbf{y}$  and  $\mathbf{y}_n$  coupled with the significant change in the genotypic correlations (particularly for yield and plant height) seem to be in conflict with each other. For instance, if the correlation between  $\mathbf{y}$  and  $\mathbf{y}_n$  is high and the  $R^2_{\text{Correlated}}$  is also large, one would expect  $\mathbf{g}$  and  $\mathbf{g}_n$  to be similar since these statistics are calculated with marker effects estimated using either  $\mathbf{y}$  or  $\mathbf{y}_n$ . However,  $\mathbf{g}$  and  $\mathbf{g}_n$  were different for trait combinations with significant correlations indicating a potential underlying difference in the genomewide marker effects.

#### *Response to Genomewide Selection and Level of Inbreeding*

Across both models, substantial gains were seen in Population 1 from genomewide selection. Compared to Cycle 0, Cycle 3 from the Control Model demonstrated a 27% gain in the selection index while Cycle 4 in the Independent Model demonstrated an 18% gain

in the selection index (Table 11). For Population 1, grain yield was improved for both models across all cycles; however, the improvement for grain moisture was inconsistent and plant height was increased (lower values desired). The differences in means of the same cycle were not significant ( $P = 0.05$ ) between the two models. Additionally, the slopes obtained from fitting the least-square means and the cycles for each model were not significantly different. These results indicated that in Population 1, the Independent Model was not superior to the Control Model. Despite the lack of a statistically significant difference, the cycle means for plant height were always numerically lower with the Independent Model than with the Control Model. As mentioned in the Materials and Methods, results on responses to selection in Population 2 will be available in fall 2017.

Inbreeding was observed across the cycles of selection in both populations. For Cycle 0, a different frequency of heterozygosity was observed in Population 2 than Population 1 because the former was created from a diallel of 12  $F_5$  lines. For Population 2, the Cycle 0 had a mean heterozygosity of 0.40 while Population 1 had a mean heterozygosity of 0.50 (Table 12). After three cycles of genomewide selection, both populations were nearing  $F_3$  status in terms of their inbreeding level (0.41-0.47 for Population 1 and 0.38-0.39 for Population 2). For both populations and across both models, the coefficient of inbreeding decreased in each successive cycle except for Cycle 2 for Population 2 using the Independent Model. In Cycle 1 of Population 2, the Independent Model led to a significantly lower ( $P = 0.05$ ) mean heterozygosity than the Control model. All other cycles of both populations did not have significantly different inbreeding levels. The Independent Model and the Control Model therefore did not alter the level of inbreeding at different rates.

## *Conclusions*

Overall, the yield trial results for Population 1 indicated a lack of a difference between the Independent Model and the Control Model. Additionally, the line estimates for each trait in the Independent Model and the Control model were highly correlated (0.93-0.99 for Population 1 and 0.97-0.99 for Population 2) and the marker effects estimated using these phenotypic values were also highly correlated between models. On the other hand, there is some evidence that the two models did not apply the same selection pressure. For example, the genotypic correlations calculated from marker effects differed between the models, particularly for the grain yield-plant height combination. As previously stated, in Population 1 the genotypic correlation between grain yield and plant height was 0.54 with the Control Model versus -0.11 with the Independent Model. Lastly, the Independent Model led to germplasm that was consistently (but not significantly) shorter at each cycle.

These results suggest that the phenotypic value of a quantitative trait may not be easily partitioned into correlated and independent portions. We speculate that the large linkage blocks that are found in biparental populations (Smith et al., 2008; Sleper and Bernardo, 2016) confound our ability to separate the genetic effects. In particular, QTL can be partitioned into three categories: independent of QTL controlling other traits (no linkage or pleiotropy); linked to other QTL that control other traits; or pleiotropic. The independent portion of a trait results only from the QTL in the first category (no linkage and no pleiotropy). However, linkage disequilibrium is likely to persist in maize breeding crosses because (i) the presence of many QTL in a genome of finite size naturally leads to linkage, and (ii) linkage disequilibrium is maximized in a cross (such as Populations 1 and 2) between two homozygous lines (Dudley, 1992). For quantitative traits in double haploid



populations derived from a biparental cross, the proportion of the genetic variance explained by independent loci is therefore likely small. As previously reported for Population 1 (Sleper and Bernardo, 2016), additional recombination may be needed to truly separate linked loci. Sleper and Bernardo (2016) demonstrated the effect of an additional meiotic event potentially contributed to disrupting the underlying repulsion and coupling blocks in a population. Unfortunately, the need for additional recombination might be a double-edged sword. Disrupting linkage blocks to estimate independent loci will also disrupt advantageous coupling blocks. Additionally, more recombination decreases the linkage disequilibrium between markers and QTL, thus reducing genomewide selection accuracies (Lian, et al., 2014).

Overall, selecting for independent portions of a trait using genomewide selection was not more effective than selecting on the entire trait in a biparental population. Once again, it appears that plant breeders are playing a game of “tug of war” between conserving valuable linkage blocks and additional recombination to induce genetic variance (Rasmussen and Phillips, 1997). Because of this phenomenon, there may not be a current, effective way to control unfavorable correlations among traits. Future explorations into targeted recombination could provide a very powerful tool to help unlock unfavorable correlations among traits (Bernardo, 2017). In the time being, utilizing a selection index to account for correlated traits in genomewide selection is recommended.

Table 9: Summary statistics for training populations used for genomewide selection in maize.

Trait	Population 1				Population 2			
	Mean	V <sub>G</sub>	V <sub>R</sub>	$h^2$	Mean	V <sub>G</sub>	V <sub>R</sub>	$h^2$
Yield (Mg ha <sup>-1</sup> )	12.72	0.49	2.25	0.57	11.05	0.78	4.53	0.55
Moisture (g kg <sup>-1</sup> )	233.0	67.3	213.2	0.66	188.8	80.9	142.8	0.80
Plant height (cm)	237.7	102.0	145.1	0.81	246.3	53.2	127.7	0.74

<sup>†</sup>All V<sub>G</sub> estimates were significant (P = 0.05).

Table 10: Pairwise phenotypic correlations between the full trait ( $y$ ), the independent portion of the trait ( $y_n$ ), and the correlated portion of the trait ( $y_R$ ). Genotypic correlations calculated using analysis of covariance ( $r_G$ ) and using molecular markers for the full trait ( $g$ ), the independent portion of the trait ( $g_n$ ), and the correlated portion of the trait ( $g_R$ ).

Trait 1	Trait 2	Population 1						
		Phenotypic			Genotypic	Marker		
		$y$	$y_n$	$y_R$	$g$	$g$	$g_n$	$g_R$
Yield	Moisture	-0.21 <sup>†</sup>	0.02	-0.25 <sup>†</sup>	-0.18 <sup>†</sup>	-0.19 <sup>†</sup>	0.02	-0.26 <sup>†</sup>
Yield	Plant Height	0.59 <sup>†</sup>	-0.10 <sup>†</sup>	0.57 <sup>†</sup>	0.32 <sup>†</sup>	0.54 <sup>†</sup>	-0.11 <sup>†</sup>	0.54 <sup>†</sup>
Moisture	Plant Height	-0.10 <sup>†</sup>	-0.11 <sup>†</sup>	-0.90 <sup>†</sup>	0.02	-0.04	-0.05	-0.92 <sup>†</sup>

Trait 1	Trait 2	Population 2						
		Phenotypic			Genotypic	Marker		
		$y$	$y_n$	$y_R$	$g$	$g$	$g_n$	$g_R$
Yield	Moisture	0.03	0.00	-0.99 <sup>†</sup>	0.06	0.05	0.00	-0.99 <sup>†</sup>
Yield	Plant Height	0.19 <sup>†</sup>	-0.14 <sup>†</sup>	0.22 <sup>†</sup>	0.31 <sup>†</sup>	0.23 <sup>†</sup>	-0.08	0.25 <sup>†</sup>
Moisture	Plant Height	-0.15 <sup>†</sup>	-0.03	-0.17 <sup>†</sup>	0.05	-0.06	0.05	-0.22 <sup>†</sup>

<sup>†</sup>Correlation was significant ( $P = 0.05$ ).

Table 11: Testcross performance of different cycles of genomewide selection for maize Population 1.

Cycle	Yield (Mg/ha)		Moisture (g/kg)		Plant height (cm)		Selection index	
	Control	Independent	Control	Independent	Control	Independent	Control	Independent
0	12.70		200.6		224.3		5.05	
1	13.27	13.37	201.8	206.8	236.8	229.8	5.34	5.52
2	13.00	13.12	202.1	201.9	236.0	224.5	5.07	5.46
3	14.16	13.05	194.0	196.0	236.0	229.8	6.39	5.36
4	14.09	13.72	192.6	201.0	243.0	230.8	6.18	5.95

†For each cycle, the testcross least-square mean of the Control Model was not significantly ( $P = 0.05$ ) different than the testcross least-square mean of the Independent Model for all three traits and the selection index

Table 12: Inbreeding level and mean heterozygosity between the Control Model and the Independent Model for both maize populations.

Cycle	Inbreeding coefficient		Mean heterozygosity	
	Control	Independent	Control	Independent
Population 1				
0	-	-	0.50	0.50
1	0.10	0.17	0.45	0.42
2	0.31	0.33	0.35	0.34
3	0.41	0.47	0.29	0.27
Population 2				
0	-	-	0.40	0.40
1 <sup>†</sup>	0.05	0.28	0.38	0.29
2	0.33	0.27	0.27	0.29
3	0.38	0.39	0.25	0.24

<sup>†</sup>Indicates a significant difference ( $P = 0.05$ ) in the inbreeding coefficient and mean heterozygosity between the Control Model and Independent Model.

## Bibliography

- Aguilar, I., I. Misztal, S. Tsuruta, G.R. Wiggans, and T.J. Lawlor. 2011. Multiple trait genomic evaluation of conception rate in Holsteins. *J. Dairy Sci.* 94:2621-2624.
- Arbelbide M., and R. Bernardo. 2004. Random mating before selfing in maize BC<sub>1</sub> populations. *Crop Sci.* 44:401-404.
- Baker, R.J. 1986. *Selection Indices in Plant Breeding*. CRC Press, Boca Raton, FL.
- Bao, Y., J.E. Kurle, G. Anderson, and N.D. Young. 2015. Association mapping and genomic prediction for resistance to sudden death syndrome in early maturing soybean germplasm. *Mol. Breeding.* 35:1-14.
- Bates, D., M. Maechler, B. Bolker, and S. Walker. 2013. lme4: Linear mixed-effects models using Eigen and S4. R package version 1.4.
- Bernardo, R. 2009. Should maize doubled haploids be induced among F<sub>1</sub> or F<sub>2</sub> plants? *Theor. Appl. Genet.* 119:255-262.
- Bernardo, R. 2010. *Breeding for quantitative traits in plants*. 2nd ed. Stemma Press, Woodbury, Minnesota.
- Bernardo, R. 2017. Prospective targeted recombination and genetic gains for quantitative traits in maize. *Plant Genome*. doi:10.3835/plantgenome2016.11.0118.
- Beyene, Y., K. Semagn, S. Mugo, A. Tarekegne, R. Babu, B. Meisel, P. Sehabiague, D. Makumbi, C. Magorokosho, S. Oikeh, J. Gakunga, M. Vargas, M. Olsen, B. M. Prasanna, M. Banziger, and J. Crossa. 2015. Genetic gains in grain yield through genomic selection in eight bi-parental maize populations under drought stress. *Crop Sci.* 55:154–163.

- Bobko, P. 2001. Correlation and regression: Applications for industrial organizational psychology and management. 2nd ed. Sage Publications, Inc., Thousand Oaks, CA.
- Calus, M.P.L., T.H.E. Meuwissen, A.P.W. de Roos, and R.F. Veerkamp. 2008. Accuracy of genomic selection using different methods to define haplotypes. *Genetics* 178:553-561.
- Calus, M.P.L., and R.F. Veerkamp. 2011. Accuracy of multi-trait genomic selection using different methods. *Genet. Sel. Evol.* 43:26.
- Chi, R.K., S.A. Eberhart, and L.H. Penny. 1969. Covariances among relatives in a maize variety (*Zea mays* L.). *Genetics* 63:511.
- Clark, R. M., T.N. Wagler, P. Quijada, and J. Doebley. 2006. A distant upstream enhancer at the maize domestication gene *tb1* has pleiotropic effects on plant and inflorescent architecture. *Nature genetics*. 38:594-597.
- Combs, E., and R. Bernardo. 2013. Genomewide selection to introgress semidwarf maize germplasm into US Corn Belt inbreds. *Crop Sci.* 53:1427-1436.
- Covarrubias-Prieto J. 1987. Genetic variability in F<sub>2</sub> maize populations before and after random mating. Ph.D. dissertation, Iowa State University, Ames, USA.
- Cuyabano, B.C.D., G. Su, and M.S. Lund. 2014. Genomic prediction of genetic merit using LD-based haplotypes in the Nordic Holstein population. *BMC Genomics* 15:1171.
- Da, Y. 2015. Multi-allelic haplotype model based on genetic partition for genomic prediction and variance component estimation using SNP markers. *BMC Genetics* 16:144.
- Daetwyler, H.D., B. Villanueva, and J.A. Woolliams. 2008. Accuracy of predicting the genetic risk of disease using a genomewide approach. *PLoS ONE* 3:e3395.

- Dudley, J.W. 1992. Theory for identification of marker locus-QTL associations in population by line crosses. *Theor. Appl. Genet.* 85:101-104.
- Endelman, J.B. 2011. Ridge regression and other kernels for genomic selection with R package rrBLUP. *The Plant Genome* 4:250-255.
- Falconer, D.S., and T.F.C. Mackay. 1996. *Introduction to Quantitative Genetics*. 4<sup>th</sup> ed. Longman, Burnt Mill, England.
- Guo, G., F. Zhao, Y. Wang, Y. Zhang, L. Du, and G. Su. 2014. Comparison of single-trait and multiple-trait genomic prediction models. *BMC genetics*. 15:1.
- Grenier, C., T. V. Cao, Y. Ospina, G. Quintero, M.H. Châtel, J. Tohme, B. Courtois, and N. Ahmadi. 2015. Accuracy of genomic selection in a rice synthetic population developed for recurrent selection breeding. *PloS ONE* 10:8.
- Haldane, J.B.S., and C.H. Waddington. 1931. Inbreeding and linkage. *Genetics* 16:358-374.
- Halauer, A.R., and M.J. Carena. 2012. Recurrent selection methods to improve germplasm in maize. *Maydica* 57:266-283.
- Hayashi, T., and H. Iwata. 2013. A Bayesian method and its variational approximation for prediction of genomic breeding values in multiple traits. *BMC bioinformatics*. 14:1.
- Hazel, L.N. 1943. The genetic basis for constructing selection indexes. *Genetics*. 28:476-490.
- Hazel, L.N., and J.L. Lush. 1942. The efficiency of three methods of selection. *Journal of Heredity*, 33(11), 393-399.
- Jacobson, A. L. Lian, S. Zhong, and R. Bernardo. 2014. General combining ability model for genomewide selection in a biparental cross. *Crop Sci.* 54.3:895-905.



- Jannink, J., and T.E. Abadie. 1999. Inbreeding method effects on genetic mean, variance, and structure of recurrent selection populations. *Crop Sci.* 39:988-997.
- Jia, Y., and J. Jannink. 2012. Multiple-Trait genomic selection methods increase genetic value prediction accuracy. *Genetics* 192:1513-1522.
- Lande, R., and R. Thompson. 1990. Efficiency of marker-assisted selection in the improvement of quantitative traits. *Genetics* 124:743-756.
- Lenth, R.V., and M. Hervé. 2015. lsmeans. R package version 2.19.
- Lemmon, Z.H., and J.F. Doebley. 2014. Genetic dissection of a genomic region with pleiotropic effects on domestication traits in maize reveals multiple linked QTL. *Genetics*. 198:345-353.
- Lian, L., A. Jacobson, S. Zhong, and R. Bernardo. 2014. Genomewide prediction accuracy within 969 maize biparental populations. *Crop Sci.* 54.4:1514-1522.
- Longin, C.F.H. 2008. Optimum allocation of test resources and comparison of alternative breeding schemes for hybrid maize breeding with doubled haploids. Doctoral dissertation, University of Hohenheim, Germany.
- Lu, H., J. Romero-Severson, and R. Bernardo. 2002. Chromosomal regions associated with segregation distortion in maize. *Theor. Appl. Genet.* 105:622-628.
- Massman, J.M., H.J.G. Jung, and R. Bernardo. 2013. Genomewide selection versus marker-assisted recurrent selection to improve grain yield and stover-quality traits for cellulosic ethanol in maize. *Crop Sci.* 53:58-66.
- Melchinger, A.E., H.H. Geiger, and F.W. Schnell. 1986. Epistasis in maize (*Zea mays* L.).  
2. Genetic effects in crosses among early flint and dent inbred lines determined by three methods. *Theor. Appl. Genet.* 72:231-239.

- Melchinger, A.E., W. Schmidt, and H.H. Geiger. 1988. Comparison of testcrosses produced from F<sub>2</sub> and first backcross populations of maize. *Crop Sci.* 28:743-749.
- Meuwissen, T.H.E., B.J. Hayes, and M.E. Goddard. 2001. Prediction of total genetic value using genome-wide dense marker maps. *Genetics* 157:1819–1829.
- Mikel, M.A., and J.W. Dudley. 2006. Evolution of North American dent corn from public to proprietary germplasm. *Crop Sci.* 46:1193-1205.
- Murigneux, A., S. Baud, and M. Beckert. 1993. Molecular and morphological evaluation of doubled-haploid lines in maize. 2. Comparison with single-seed-descent lines. *Theor. Appl. Genet.* 87:278-287.
- Ooijen, J.W., and R.E. Voorrips. 2002. JoinMap: version 3.0: software for the calculation of genetic linkage maps. Wageningen University and Research Center, The Netherlands.
- R Development Core Team. 2012. R: a language and environment for statistical computing. Vienna: R Foundation for Statistical Computing. <http://www.r-project.org>. Accessed Dec 2015.
- Rasmusson, D.C., and R.L. Phillips. 1997. Plant breeding progress and genetic diversity from de novo variation and elevated epistasis. *Crop Sci.* 37:303-310.
- Riggs, T.J., and J.W. Snape. 1977. Effects of linkage and interaction in a comparison of theoretical populations derived by diploidized haploid and single seed descent methods. *Theor. Appl. Genet.* 49:111-115.
- Ross, A.J. 2002. Genetic analysis of ear length and correlated traits in maize. Ph.D. dissertation, Iowa State University, Ames, USA.

- Rober, F.K., G.A. Gordillo, and H.H. Geiger HH. 2005. In vivo haploid induction in maize-performance of new inducers and significance of doubled haploid lines in hybrid breeding. *Maydica* 50:275-283.
- Rutkoski, J., J. Benson, Y. Jia, G. Brown-Guedira, J.L. Jannink, and M. Sorrells. 2012. Evaluation of genomic prediction methods for *Fusarium* head blight resistance in wheat. *Plant Gen.* 5:51-61.
- Schulthess, A.W., Y. Wang, T. Miedaner, P. Wilde, J.C. Reif, and Y. Zhao. 2016. Multiple-trait-and selection indices-genomic predictions for grain yield and protein content in rye for feeding purposes. *Theor. Appl. Genet.* 129:273-287.
- Senior, M.L., E.C.L Chin, M. Lee, J.S.C Smith, and C.W. Stuber. 1996. Simple sequence repeat markers developed from maize sequences found in the GEN-BANK database: Map construction. *Crop Sci.* 36.6:1676-1683.
- Silva, J.C., and A.R. Hallauer. 1975. Estimation of epistatic variance in Iowa Stiff Stalk Synthetic maize. *J. Hered* 66:290-296.
- Sleper, J.A., and R. Bernardo. 2016. Recombination and genetic variance among maize doubled haploids induced from F<sub>1</sub> and F<sub>2</sub> plants. *Theor. Appl. Genet.* 1-8.
- Smith, H.F. 1936. A discriminant function for plant selection. *Annals of Eugenics.* 7:240-250.
- Smith, J.S.C., T. Hussain, E.S. Jones, G. Graham, D. Podlich, S. Wall, and M. Willimans. 2008. Use of doubled haploids in maize breeding: implications for intellectual property protection and genetic diversity in hybrid crops. *Mol. Breeding* 22:51-59.
- Sokal, R.R., and C.D. Michener. 1958. A statistical method for evaluating systematic relationships. *Univ. Kansas Sci. Bull.* 38:1409-1438.

- Stubber, C.W., and R.H. Moll. 1971. Epistasis in maize (*Zea mays* L.). II. Comparison of selected with unselected populations. *Genetics* 67:137-149.
- Weir, B.S., C.C. Cockerham, and J. Reynolds. 1980. The effects of linkage and linkage disequilibrium on the covariances of noninbred relatives. *Heredity* 45:351-359.
- Ziyomo, C., and R. Bernardo. 2012. Drought tolerance in maize: indirect selection through secondary traits versus genomewide selection. *Crop Sci.* 53:1269-1275.